

HIGH-DIMENSIONAL PROBLEMS IN STATISTICS AND PROBABILITY:  
CORRELATION MINING AND DISTRIBUTED LOAD BALANCING

Miheer Dewaskar

A dissertation submitted to the faculty of the University of North Carolina at Chapel Hill  
in partial fulfillment of the requirements for the degree of Doctor of Philosophy in the  
Department of Statistics and Operations Research.

Chapel Hill  
2021

Approved by:

Sayan Banerjee

Shankar Bhamidi

Amarjit Budhiraja

Jan Hannig

Andrew B. Nobel

©2021  
Miheer Dewaskar  
ALL RIGHTS RESERVED

## ABSTRACT

Miheer Dewaskar : High-dimensional problems in statistics and probability: correlation mining  
and distributed load balancing  
(Under the direction of Shankar Bhamidi, Amarjit Budhiraja, and Andrew B. Nobel)

Technological progress has encouraged the study of various high-dimensional systems through the lens of statistics and probability. In this dissertation, we consider two such high-dimensional problems: the first arising in the integration of genomic data, and the second arising in probabilistic models for load balancing. A brief description follows.

It is now common across many scientific and engineering disciplines to have multiple types of features measured on the same set of samples. In the first part of this dissertation, in the context of two measurement types, we focus on the exploratory problem of finding bimodules: these are sets of features from the two data types that have significant aggregate cross-correlation. Based on the iterative-testing framework that has been recently used in other settings, we design a new methodology to find bimodules. We apply this methodology to the problem of eQTL-analysis in genomics to identify gene-SNP association networks.

In the second part of this dissertation, motivated by load balancing problems in large data centers, we study a processing system with multiple queues known as the Supermarket model. In this system, each incoming job is routed into one of the  $n$  available queues based on the following randomized scheme:  $d$  out of the  $n$  queues are sampled at random and the job is assigned into the smallest of the  $d$  sampled queues. Hence, when  $d = 1$ , each job joins a random queue, while for  $d = n$ , each job joins the shortest of all  $n$  queues. Here we prove functional central limit theorems for this system in various regimes where  $d$  and  $n$  scale to infinity and the system load approaches criticality.

*To my grandparents.*

## ACKNOWLEDGEMENTS

A lot of people have directly and indirectly contributed to the creation of this dissertation. In the direct sense, the work in the first part of this dissertation was conducted in collaboration with John Palowitch, Mark He, Michel I. Love, and Andrew B. Nobel. The work in the second part of this dissertation was conducted in collaboration with Shankar Bhamidi and Amarjit Budhiraja. I would also like to acknowledge grants NIH R01 HG009125-01, NSF DMS-1613072, and that from SAMSI, which allowed me to focus my time on research. I would also like to thank Fred Wright for recommending the fast p-value approximation used in Chapter 3, and for taking time out to explain its details to me. In the rest of this section, I will now acknowledge the many people who have contributed to this dissertation by helping me complete my PhD journey over the last five years.

First and foremost, I would like to thank my three amazing advisors who put in a lot of effort into mentoring and supporting me. They continue to inspire me to work hard and to improve my abilities. I will start by thanking Andrew for all the confidence that he has shown in me. Were it not for his support and training, I might not have pursued research in statistics. I am grateful that he kept me regularly accountable and provided me with a challenging environment that helped to improve my work. There were times in the research process where I did not know what I was doing, but I am glad that Andrew continued to nudge and support me towards completing the work. Apart from the several technical skills that I admire in him, I would like emulate his energy, work ethic, and interpersonal skills.

Next, I would like to thank Amarjit for being, throughout the years, a kind and patient voice of reason for me. He has always been able to provide satisfying answers to the numerous questions on math, research, career, and life that I have posed to him. His nuggets of wisdom have had an important influence on me. Apart from the several interesting mathematical techniques that I learned from Amarjit, by observing him work and teach, I have also learned about the importance of patience, discipline, and humility, while doing mathematics.

Finally, I would like to thank Shankar for being an amazing mentor and a friend. His curiosity and passion towards any task, be it research, teaching, and even administrative duties, has been inspiring to watch. I appreciate how he thinks about various ways to help everyone around him, particularly his students. His hard work and energy is contagious, and I aspire to emulate the numerous positive qualities in him. Much of the research and career advice he has given me over the years has proved invaluable: be it the various books that he suggested that I read, or the importance of focusing on the process over the outcome when doing creative work.

I would also like to thank my other committee members, Sayan Banarjee and Jan Hannig, for finding time to provide pertinent feedback on my work, which has helped me improve its presentation. Additionally, over the years, Sayan's big-brotherly advice on the research process, job search, career advice etc. has been quite reassuring.

Next, I would like to thank the STOR department and all its people who created an intellectually stimulating and nurturing environment. This includes Professors Kulkarni, Ziya, Lu, Tran-Dihn, and Pataki, whose classes I really enjoyed during my first two years. I thank Prof. Kulkarni for his support during my initial years, including the flexibility in choosing my course work, for offering me Tilgul during Sankranti each year, and for organizing his yearly talent parties.

I would also like to thank the departmental staff who would promptly resolve any problems that I would bring up to them. Particularly, I would like to thank Christine Keat for sending reference letters to various jobs that I applied to, and for promptly assisting me with various questions and issues that arose as an international student.

Next, I would like to thank the fellow students from the department that have played an important role in helping me thrive in the past five years. This includes

- Awesome seniors like John Palowitch, Eric Friedlander, Jonathan Williams, Jimmy Jin, and Duyeol Lee, who provided me with help and advice throughout the years.
- Colleagues like Kevin O'Connor, Brendan Brown, Weibin Mo, Jack Prothero, Mark He, Gang Li, Jonghwan Yoo, and Kentaro Hoffman with whom it was fun to discuss and debate big and small ideas, and
- Friendships with Carson Mosso, Michel Bostwick, Wei Liu, and Anand Bhatia that I will always cherish.

Finally, I would like to thank my friends that are almost family – Ajeenchya, Amrita, Aman, Samo, Benjy, and Manish. Without them I can not imagine surviving my life outside the department in the past five years. I would also like to thank Seema tai’s family for considering me as one of their own and providing me a home in NC. Finally, I would like to thank my family for the constant support I received from them. This includes Aaji and Aajoba, who stayed with me for the first six months, so that I could settle in and focus on my work. Mama and Mami, who always try to help and guide me in various ways that they can, and finally Aai and Baba who are always there for me.

## TABLE OF CONTENTS

LIST OF TABLES .....	x
LIST OF FIGURES.....	xi
LIST OF ABBREVIATIONS .....	xii
LIST OF SYMBOLS.....	xiii
1 Introduction .....	1
1.1 Detecting bimodules: sets of cross-correlated features in multi-view data.....	2
1.2 Limit theorems for the Supermarket model with growing choices .....	6
<b>Correlation Mining .....</b>	<b>11</b>
2 Literature review for bimodule detection .....	11
2.1 Introduction.....	11
2.2 Survey of methods for bimodule detection.....	12
3 Bimodule Search Procedure .....	16
3.1 Notation and stochastic setting .....	16
3.2 Stable population bimodules .....	17
3.3 Sample stable bimodules and the Bimodule Search Procedure (BSP) .....	21
3.4 BSP implementation details .....	26
3.5 The continuous permutation scheme .....	33
4 Data analysis .....	47
4.1 Simulation study .....	47
4.2 Real data study: bimodules for eQTL analysis .....	56
4.3 Details of data analysis .....	62



<b>Distributed load balancing</b> .....	<b>68</b>
5 Literature review for distributed load balancing .....	68
5.1 The balls and bins problem .....	68
5.2 The Supermarket model .....	70
6 Maximum load in the balls and bins problem .....	75
6.1 Model description and overview .....	75
6.2 Concentration and its consequences .....	77
6.3 Fluid limit approximation .....	82
6.4 Estimating the expected maximum .....	90
6.5 Technical estimates .....	92
7 Limit theorems for the Supermarket model .....	95
7.1 Introduction .....	95
7.2 Main Results .....	99
7.3 Poisson Representation of State Processes .....	110
7.4 The Law of Large Numbers .....	113
7.5 Properties of the Near Fixed Point .....	122
7.6 Preliminary estimates under diffusion scaling .....	124
7.7 Proof of Theorem 9 .....	131
7.8 Proof of Theorem 10 .....	141
7.9 Proof of Theorem 11 .....	148
7.10 Technical estimates .....	159
8 Conclusion .....	167
8.1 Correlation mining : summary and future directions .....	167
8.2 Distributed load balancing : summary and future directions .....	169
A Gene Ontology enrichment of bimodules .....	171
BIBLIOGRAPHY .....	179

## LIST OF TABLES

Table 1.1	Regimes for analysis of the Supermarket model .....	7
Table 4.1	Comparison of BSP and standard eQTL analysis.....	61

## LIST OF FIGURES

Figure 1.1	Multi-view data and bimodules .....	3
Figure 4.1	Recovery of true bimodules in the simulation study .....	51
Figure 4.2	Recovery of true bimodules on increasing simulation sample size.....	53
Figure 4.3	Bimodule sizes under real and permuted data .....	59
Figure 4.4	Correlations of essential-edges from BSP bimodules .....	59
Figure 4.5	Genomic plots of two BSP bimodules .....	62
Figure 4.6	Sizes of sCCA bimodules .....	63
Figure 4.7	Network statistics for BSP bimodules .....	66
Figure 4.8	Connectivity of BSP bimodules under standard eQTL analysis .....	66
Figure 4.9	Genomic plots for nine more BSP bimodules .....	67
Figure 5.1	The Supermarket model with parameters $(n, d_n, \lambda_n)$ .....	71

## LIST OF ABBREVIATIONS

BLAS	Basic Linear Algebra Subprograms
BSP	Bimodule Search Procedure
B.Y.	The Benjamini–Yekutieli multiple testing procedure
CCA	Canonical Correlation Analysis
CDF	Cumulative Distribution Function
CLT	Central Limit Theorem
CONDOR	Method by Platig et al. (107) to find communities in eQTL networks
CTMC	Continuous Time Markov Chain
eQTL	Expression Quantitative Trait Loci
FLLN	Functional Law of Large Numbers
FCLT	Functional Central Limit Theorem
GO	The Gene Ontology database (50)
JSQ	Join the Shortest Queue routing scheme
LLN	Law of Large Numbers
OU	Ornstein–Uhlenbeck process
PEER	Probabilistic Estimation of Expression Residuals
RCLL	Right Continuous with Left Limits
SBM	Standard Brownian Motion
sCCA	Sparse Canonical Correlation analysis
SNP	Single Nucleotide Polymorphism

## LIST OF SYMBOLS

$\mathbb{N}, \mathbb{N}_0$	The set of numbers $\{1, 2, \dots\}$ and $\{0, 1, 2, \dots\}$ , respectively
$[n]$	The set $\{1, \dots, n\}$
$\mathbb{R}$	The real number line
$\wedge, \vee$	Minimum and maximum operations, respectively
$f(n) \gg g(n)$	This means that $\lim_n f(n)/g(n) = \infty$
$f(n) \sim g(n)$	This means that $\lim_n f(n)/g(n) = 1$
$f(n) = O(g(n))$	This means that $\limsup_n f(n)/g(n) < \infty$
$f(n) = o(g(n))$	This is the same as $f(n) \ll g(n)$
$\mathcal{N}(\mu, \sigma^2)$	Normal distribution with mean $\mu$ and variance $\sigma^2 > 0$
$\mathcal{N}_d(\mu, \Sigma)$	Multivariate Normal with mean $\mu \in \mathbb{R}^d$ and co-variance $\Sigma \in \mathbb{R}^{d \times d}$
$\mathcal{O}_k$	The space of $k \times k$ orthonormal matrices
$L \odot M$	The Hadamard product of matrices $L$ and $M$
$l_1^\downarrow$	The space of infinite non-increasing summable sequences
$l_p$	The space of infinite sequences with finite $p$ -norm
$\mathcal{S}$	A Polish space (i.e. a complete separable metric space)
$C([0, \infty) : \mathcal{S})$	The space of continuous functions from $[0, \infty)$ to $\mathcal{S}$
$\mathbb{D}([0, \infty) : \mathcal{S})$	The space of RCLL functions $[0, \infty) \rightarrow \mathcal{S}$ with Skorokhod topology
$\mathbf{E}, \mathbf{P}$	The expectation and probability operators
$\lambda_n$	Average arrival rate for the Supermarket model with $n$ servers
$G_{n,i}(t)$	The fraction of queues that have at least $i \in \mathbb{N}_0$ customers at time $t$
$\mathbf{G}_n$	The stochastic process $\mathbf{G}_n = (G_{n,i})_{i \in \mathbb{N}}$ that takes values in $l_1^\downarrow$
$\boldsymbol{\mu}_n$	The Near-Equilibrium point for the Supermarket model with $n$ servers
$\mathbf{Z}_n$	The scaled and centered process $\mathbf{Z}_n = \sqrt{n}(\mathbf{G}_n - \boldsymbol{\mu}_n)$
$\mathbf{e}_k$	The vector $(0, \dots, 0, 1, 0, \dots)$ with the $k$ th coordinate equal to one
$\mathbf{f}_k$	The vector $(1, \dots, 1, 0, \dots)$ with the first $k$ coordinates equal to one
$(\Gamma_\alpha, \hat{\Gamma}_\alpha)$	The Skorokhod map on $\mathbb{R}$ with a downward reflection at $\alpha \in (-\infty, \infty]$

## CHAPTER 1

### **Introduction**

The fields of mathematics and statistics have been unsung heroes of our scientific and technological progress. The language and machinery developed in mathematics has provided us with rigorous tools to describe and understand models of complex systems. In parallel, the field of statistics has provided us with tools to deal with uncertainties in the collection and processing of real world data. Tools from mathematics and statistics used in areas like computer science and signal processing have been crucial for the development of modern technologies. But conversely, the field of statistics in particular has also received a boost due to the advent of high-speed computation and technologies that collect and store large amounts of data. In the future, by combining modern mathematical, statistical, and technological tools, partially observed and increasingly complex real world systems could be understood on the basis of large amounts of data. Such progress certainly seems necessary if we wish to fully understand functioning of complex systems like biological processes that govern our body or micro-economic processes that govern the macro-economy. Motivated by this goal of inferring properties of increasingly complex systems, in this dissertation the author presents two topics that originate in statistics and applied probability.

### **Correlation mining and randomized load balancing**

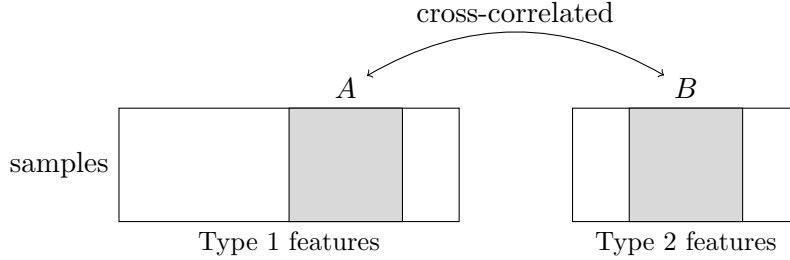
The first topic of this dissertation is related to the field of correlation mining, which aims to understand the true covariance structure among a large collection of features, based on joint observations from typically a limited number of samples. Correlation mining has applications to various exploratory data analysis tasks in science and engineering, where a large number of features may be measured for each sample, and the discovery of certain kinds of relationship between features is of interest. In fact, due to the emergence of new measurement technologies, features of several types may be measured for the same samples. Such data is sometimes called multi-view or multi-modal data, and the statistical methods that “integrate” such data provide new insights by a

combined analysis of all the data types, beyond what may be obtained by analyzing each data type separately. Along these lines, in the first part of this dissertation, we study a new method for exploratory data analysis that identifies groups of features with significant aggregate correlation between two data-types. Such a group of features, called a bimodule, represents coordinated activity between features from the two data types, and may point to shared latent variables, causal interactions, or other functional relationship between the data types. Along with the motivating example of eQTL analysis in genomics, we describe this problem further in Section 1.1.

In the second part of this dissertation, we study a processing system with multiple queue motivated by the large scale load-balancing problem in modern data centers. The simplest version this problem can be described in the setup of balls and bins: suppose  $n$  balls are to inserted into  $n$  bins, one at a time. Consider a routing scheme that sends each ball to the least loaded bin among a subset of  $d \in \{1, \dots, n\}$  bins chosen randomly for each ball. When  $d = n$  this means that each ball is put in the least loaded of all bins, while when  $d = 1$ , each ball is assigned to a randomly chosen bin. As a measure of effectiveness of our load balancing scheme, we could study the maximum bin load when all the  $n$  balls are inserted. (Note that this system is exactly balanced when the maximum is 1. This happens for instance when  $d = n$ .) Although the exact distribution of the maximum for  $d < n$  is not known, asymptotically as  $n \rightarrow \infty$ , the maximum can be shown to be  $O(\frac{\log n}{\log \log n})$  when  $d = 1$  and  $\log_d \log n + O(1)$  for a fixed  $d \geq 2$ . This drastic decrease in the the size of the maximum between  $d = 1$  (random assignment) and fixed  $d \geq 2$  (random assignment with limited choice) is called as the “power of choice”. This phenomenon has application to many areas in computer science like load-balancing, hashing, and collision protocols (112). In this dissertation we will consider the asymptotic behavior of the power-of- $d$ -choices load balancing scheme as  $d = d_n$  is now allowed to grow with  $n$ . This problem and its motivation, particularly in the continuous time setup called the Supermarket model, are further described in Section 1.2.

## 1.1 Detecting bimodules: sets of cross-correlated features in multi-view data

With the ongoing development and application of moderate- to high-throughput measurement technologies in fields such as genomics, neuroscience, ecology, and atmospheric science, researchers are often faced with the task of analyzing and comparing two or more data sets derived from



**Figure 1.1:** Illustration of multi-view data and a bimodule  $(A, B)$  (shaded). The two data matrices, measuring features of Type 1 and Type 2, are matched by samples. The bimodule has the property that the aggregate cross-correlation between features from  $A$  and  $B$  is statistically significant.

a common set of samples. In most cases, different technologies measure different features, and capture different information about the samples at hand. While one may analyze the data arising from different technologies separately, additional and potentially important insights can often be gained from the joint (or integrated) analysis of the data sets. Joint analysis, also called multi-view or multi-modal analysis, has received considerable attention in the literature, see Lahat et al. (68), Meng et al. (88), Tini et al. (120), Pucher et al. (109), McCabe et al. (83), Sankaran and Holmes (113) and the references therein for more details.

We will refer to the measurements arising from a particular technology as a data type, and will restrict our attention to problems in which two data types, referred to as Type 1 and Type 2, are under study. Our primary interest is in exploring associations between the measured features of the two data types. In particular, we wish to identify pairs  $(A, B)$ , where  $A$  is a set of features of Type 1 and  $B$  is a set of features of Type 2, such that the aggregate correlation between the features in  $A$  and those in  $B$  is large. This setup is illustrated in Figure 1.1. Throughout we consider the unsupervised setting, in which the analysis does not make use of an external response. The problem of identifying sets of highly correlated features within a single data type has been widely studied, typically through clustering and related methods. Borrowing from the use in genomics of the term “module” to refer to a set of correlated genes, we refer to the feature set pairs  $(A, B)$  of interest to us as *bimodules*. The term bimodule has also appeared, with somewhat different meaning, in (131), (103), and (99).

We will refer to the correlations between features from different data types as *cross-correlations*, and note that this usage differs from that in time-series analysis. Correlations among features of the same data type will be referred to as *intra-correlations*. Cross-correlations provide information



about interactions between features from the two data types. These interactions are of interest in many applications, for example, in studying the relationship between the characteristics of species and those of their environment (see, e.g., 42) in ecology, identifying brain regions associated with different behaviors (85) in neuroscience, and studying the relationship between temperature and precipitation in climate science (41). A bimodule provides evidence for the coordinated activity of features from different data types. Coordination may arise, for example, from common function or functional relationships, or causal interactions. Bimodules can identify potentially informative downstream analyses, or suggest the more targeted acquisition and analysis of new data. Importantly, bimodules can capture aggregate behavior, which may be significant even when no individual pair of features has high cross-correlation. As such, the search for bimodules can effectively leverage low-level signals across multiple features.

### 1.1.1 Bimodule Search Procedure

In this dissertation, we propose and analyze a method called the Bimodule Search Procedure (BSP) for identifying bimodules in moderate to high dimensional data sets. BSP is not based on formulating or fitting a statistical model, or on detailed distributional assumptions. Instead, the method relies on general multiple testing principles.

A key feature of BSP is that it seeks stable bimodules. A bimodule  $(A, B)$  is stable if  $A$  coincides with the features that are significantly associated in aggregate with the features in  $B$ , and vice versa. We examine stable bimodules in the population and sample settings. In the population setting, stability has connections with Nash equilibria (93) in a simple two-player game, and with the connectivity of the bipartite graph representing the cross-correlations of the Type 1 and Type 2 variables. The latter connection with bipartite networks is pursued throughout the first part of this dissertation, and in particular, provides a principled way to extract an association network from a bimodule.

BSP is based on an iterative testing framework that has found application in other exploratory problems, see (98, 17, 18). In the present setting, we employ fast, moment-based approximations to permutation p-values for sums of squared correlations. We emphasize that these p-values explicitly account for the intra-correlations of the features in  $A$  and  $B$ , attenuating significance when these

intra-correlations are high. BSP also has a fast implementation as an R package available at <https://github.com/miheerdew/cbce>.

### 1.1.2 Expression Quantitative Trait Loci Analysis

Much of the existing work on bimodule discovery is focused on the integrated analysis of genomic data. To motivate and provide context for this work, and the methods introduced in the first part of this dissertation, we briefly discuss the problem of expression quantitative trait loci (eQTL) analysis in genomics. An application of BSP to eQTL analysis is given in Chapter 4.

Genetic variation within a population is commonly studied by considering single nucleotide polymorphisms, called SNPs. A SNP is a single base pair site in the genome where there is allelic variation in the population. The dosage of the SNP for an individual in terms of one of the alleles can be considered taking values 0, 1, or 2, which we will refer to as the “value” at the SNP. After normalization and covariate correction, the value of a SNP may no longer be discrete.

eQTL analysis seeks to identify SNPs that affect the expression of one or more genes; a SNP-gene pair for which the expression of the gene is correlated with the value of the SNP is referred to as an eQTL. Identification of eQTLs is an important first step in the study of genomic pathways and networks that underlie disease and development in human and other populations (94, 1).

In modern eQTL studies it is common to have measurements of 10-20 thousand genes and 2-5 million SNPs on hundreds (or in some cases thousands) of samples. Identification of putative eQTLs or genomic “hot spots” is carried out by evaluating the correlation of numerous SNP-gene pairs, and identifying those meeting an appropriate multiple testing based threshold. In studies with larger sample sizes it may be feasible to carry out *trans*-eQTL analyses, which consider all SNP-gene pairs regardless of genomic location. However, it is more common to carry out *cis*-eQTL analyses, in which one restricts attention to SNP-gene pairs for which the SNP is within some fixed genomic distance (often 1 million base pairs) of the gene’s transcription start site, and in particular, on the same chromosome (c.f. 126, 54). We use the prefixes *cis*- and *trans*- to refer to the type of eQTL analysis, while using adjectives *local* and *distal* to denote the proximity of the discovered SNP-gene pairs. In particular, although *cis*-eQTL analysis focuses on finding local eQTLs, *trans*-eQTL analysis can discover *both* local and distal eQTLs.

As a result of multiple testing correction needed to address the large number of SNP-gene pairs under study, both *trans*- and *cis*-eQTL analyses can suffer from low power. Several methods have been proposed to improve the power of standard eQTL analysis, including penalized regression schemes that try to account for intra-gene or intra-SNP interaction networks (119, and references therein) and methods that consider gene modules as high-level phenotypes to reduce the burden of multiple-testing (65). Alternatively, one may also shift attention from individual SNP-gene pairs to SNP-gene bimodules, that is, to sets of SNPs and sets of genes having large aggregate cross-correlation. A number of bimodule search methods have been proposed and developed in the context of eQTL analysis. These include methods based on Gaussian graphical models (34, 32), bipartite community detection (107, 47), penalized regression (30), and sparse canonical correlation analysis (sCCA) (101, 102). These methods have the potential to enhance and improve standard eQTL approaches that focus primarily on identifying significant SNP-gene pairs. As SNPs and genes often act in concert with one another, bimodule discovery methods can gain statistical power from group-wise interactions by borrowing strength across individual SNP-gene pairs.

### 1.1.3 Organization of the first part of this dissertation

We will provide a literature review on methods from correlation mining and multi-view data analysis that are relevant to bimodule detection in Chapter 2. Next, in Chapter 3, we will describe the notion of stable bimodules and the Bimodule Search Procedure. Then, in Chapter 4, we will see analysis from BSP and the competing methods CONDOR and sCCA on a complex simulation study and on the latest version (v8) of the Thyroid data, which was obtained from the GTEx consortium for eQTL analysis.

## 1.2 Limit theorems for the Supermarket model with growing choices

In data centers with a large number of servers, incoming requests have to be quickly routed to a server while simultaneously making best use of the computational resources of the available servers. One of the simplest mathematical models to analyze this problem is the *Supermarket model*: there are  $n \in \mathbb{N}$  servers, each with their own queue (like in a supermarket checkout line), independently processing jobs at rate 1. Requests arrive to the entire system at rate  $n\lambda_n$ , and

Reference	Regimes of $a$ , $b_n$ and $k$	Analysis type
Braverman (23)	$a = 0.5$ , $b_n = 1$ , $k = 1$	Convergence of stationary distribution
Eschenfeldt and Gamarnik (45)	$a = 0.5$ , $b_n = 1$ , $k = 1$	Functional central limit theorem
Mukherjee et al. (91)	$a = 0.5$ , $b_n \in ((0.5 + \frac{\log \log n}{\log n}, 1]$ , $k = 1$	Functional central limit theorem
Theorem 11 ( $\alpha = 0$ )	$a \in (1/2, 1)$ , $b_n \in [a + \frac{\log \log n}{\log n}, 1]$ , $k = 1$	Functional central limit theorem
Theorem 11 ( $\alpha \in (0, \infty)$ )	$a = 0.5$ , $b_n \in [0.5 + \frac{\log \log n}{\log n}, 1]$ , $k = 1$	Functional central limit theorem
Theorem 11 ( $\alpha = \infty$ )	$a \in [1/3, 1/2)$ , $b_n \in ((0.5, 1]$ , $k = 1$	Functional central limit theorem
Theorem 10 ( $\alpha = \infty$ )	$a \in (1/4, 1/2)$ , $b_n = 0.5$ , $k = 1$	Functional central limit theorem
Theorem 9 ( $\alpha = 0$ )	$a \in (0, 1)$ , $b_n = (a + \frac{\log \log n}{\log n})/k \rightarrow b$ , $2a < 1 + b(k - 1)$ , $a/b \in \mathbb{N}$	Functional central limit theorem
Brightwell et al. (24)	$a \in (0, 1)$ , $b_n \rightarrow b \in (0, 1]$ , $2a < 1 + b(k - 1)$ , $a/b \notin \mathbb{N}$	Equilibrium queue lengths
Liu and Ying (74)	$a \in (0, 1/2)$ , $b_n \in [a + \frac{\log n \log n}{\log n}, 1]$ , $k = 1$	Equilibrium performance

**Table 1.1:** Existing results and asymptotic regimes for the Supermarket model covered by theorems in Chapter 7. Here  $1 - \lambda_n = n^{-a}$ ,  $d_n = \lceil n^{b_n} \rceil$ , and  $k \doteq \lim_n \lceil a/b_n \rceil$  (the average time spent in the system). The notation  $b_n \in ((l_n, r])$  denotes the condition  $n^{l_n} \ll d_n \leq n^r$ .

need to be instantaneously routed to the waiting queue of some server. We will assume throughout that we are working with the system in the Markovian setup, i.e. inter-arrival and service times are mutually independent and exponentially distributed. While various routing schemes can be used for the Supermarket model (121), we consider the *power-of- $d$*  scheme described earlier: each incoming request examines  $d$  randomly chosen servers without replacement, and joins the shortest queues among the  $d$  sampled servers.

Although much of the previous work on the theoretical analysis of the Supermarket model with the power-of- $d$  routing scheme surrounded understanding the system for a fixed  $d$ , recent work (91, 45, 24, 73, 74) has allowed  $d = d_n$  to increase with  $n$  while working in the heavy traffic regime ( $\lambda_n \uparrow 1$ ). Interests along these lines arose when Eschenfeldt and Gamarnik (45), via a diffusion limit result, showed that  $d_n = n$  achieves optimal load balancing (performing similar in terms average delay to when a single global queue is maintained) as  $n \rightarrow \infty$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta > 0$ . Although no fixed value of  $d$  can achieve such a performance (121), by using coupling arguments, it is shown in Mukherjee et al. (91) that the diffusion limit (and hence the optimal performance) of  $d_n = n$  continues to hold as long as  $d_n \gg \sqrt{n} \log n$ .

Motivated by the above results, we will study the asymptotic behavior of the Supermarket model for various growth rates of  $d_n \rightarrow \infty$ . Table 1.1 provides an overview of the diffusion limit theorems in Chapter 7, and how they compare to previous work. The following is an overview of the main results in Chapter 7.

The Supermarket model with  $n$  servers can be described by the infinite dimensional stochastic process  $\mathbf{G}_n(t) = (G_{n,1}(t), G_{n,2}(t), \dots)$ , where  $G_{n,i}(t)$  for any  $i \geq 1$  is the fraction of queues at time  $t$

that have at least  $i$  customers, including the one in service. This process takes values in the space  $l_1^\downarrow$  of non-increasing summable sequences. One can show the following semi-martingale representation for  $\mathbf{G}_n$  obtained from time changed Poisson processes (details in Section 7.3):

$$\mathbf{G}_n(t) = \mathbf{G}_n(0) + \int_0^t \{\mathbf{a}_n(\mathbf{G}_n(s)) - \mathbf{b}(\mathbf{G}_n(s))\} ds + \mathbf{M}_n(t), \quad (1.1)$$

where  $\mathbf{a}_n$  and  $\mathbf{b}$  are two functions from  $l_1^\downarrow$  to  $l_1$ ,  $\mathbf{b}$  is linear, and  $\mathbf{M}_n$  is a martingale.

### Functional law of large numbers

First, in Theorem 8, we will prove that  $\mathbf{G}_n$  converges to a deterministic limit as  $n \rightarrow \infty$ ,  $d_n \rightarrow \infty$  and  $\lambda_n \rightarrow \lambda \in [0, \infty)$ . For this we will show that  $\mathbf{M}_n$  converges to zero,  $\mathbf{G}_n$  is tight in a suitable functional space, and that any limit point of  $\mathbf{G}_n$  satisfies a collection of constrained differential equations governed by Skorokhod maps that emerge in (1.1) from the limiting behavior of  $\mathbf{a}_n$ . The limiting differential equations do not depend on the growth-rate of  $d_n$  and are equivalent to those identified in (91). However our formulation in terms of Skorokhod maps will allow us to deal with non-differentiability and prove uniqueness of the limiting system, completing the program for the law of large numbers started in (91).

### Functional central limit theorems around the near-equilibrium point

Next, we will prove central limit theorems (CLTs) as  $n \rightarrow \infty$ ,  $d_n \rightarrow \infty$  and  $\lambda_n \uparrow 1$  for  $\mathbf{Z}_n \doteq \sqrt{n}(\mathbf{G}_n - \boldsymbol{\mu}_n)$ , where  $\boldsymbol{\mu}_n \in l_1^\downarrow$  is the unique solution to  $\mathbf{a}_n(\boldsymbol{\mu}_n) = \mathbf{b}(\boldsymbol{\mu}_n)$ . Here  $\boldsymbol{\mu}_n$  is called as the *near-equilibrium point* because it is the unique state of the system where the inflow and outflow rates for each coordinate become equal. Also, notice that the drift term in (1.1) vanishes at  $\boldsymbol{\mu}_n$  and, as previously noted,  $\mathbf{M}_n$  is converging to zero. Hence if  $\mathbf{G}_n(0) = \boldsymbol{\mu}_n$  then  $\mathbf{G}_n$  will tend to remain close to  $\boldsymbol{\mu}_n$  over finite intervals of time. For the CLT, we will center the  $\mathbf{G}_n$  process around  $\boldsymbol{\mu}_n$ , while for the CLT proved in (45, 91) when  $d_n \gg \sqrt{n} \log n$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta > 0$ , centering is considered around  $\mathbf{e}_1 = (1, 0, \dots)$ . The reason for a different choice of centering is that, as shown in (91), the process  $\sqrt{n}(\mathbf{G}_n - \mathbf{e}_1)$  is not tight when  $d_n \ll \sqrt{n} \log n$ .

We will reconcile the above differences by showing that  $\mathbf{e}_1 = (1, 0, \dots)$  is a fixed point of the fluid limit (i.e. functional law of large numbers) obtained in Theorem 7. However, the fluid limit has

uncountably many fixed points given by  $\mathbf{f}_k^\gamma = \sum_{i=1}^k \mathbf{e}_i + \gamma \mathbf{e}_{k+1} \in l_1^k$  for  $k \in \mathbb{N}$  and  $\gamma \in [0, 1)$ , where  $\mathbf{e}_i \in l_1$  is the unit vector whose  $i$ th coordinate is non-zero. Interestingly, the near-equilibrium point  $\boldsymbol{\mu}_n$  can converge to all of these fixed points for various choices of  $\lambda_n \uparrow 1$  and  $d_n \rightarrow \infty$ . Hence our results described in Chapter 7 will recover those in (45, 91) for  $\sqrt{n}(G_n - \mathbf{e}_1)$  when  $d_n \gg \sqrt{n} \log n$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta > 0$  along with the interesting edge cases of  $\beta = 0$  and  $d_n \sim \sqrt{n} \log n$  under suitable conditions (see Corollary 6, Chapter 7). In addition, we will also obtain limit theorems for  $\mathbf{Z}_n$  in the three canonical regimes: (a)  $d_n/\sqrt{n} \rightarrow 0$ , (b)  $d_n/\sqrt{n} \rightarrow c \in (0, \infty)$ , and (c)  $d_n/\sqrt{n} \rightarrow \infty$ .

Our limiting result for  $\mathbf{Z}_n$  will start from its semi-martingale representation

$$\mathbf{Z}_n(t) = \mathbf{Z}_n(0) + \int_0^t \{\mathbf{A}_n(\mathbf{Z}_n(s)) - \mathbf{b}(\mathbf{Z}_n(s))\} ds + \sqrt{n} \mathbf{M}_n(t) \quad (1.2)$$

derived from (1.1), where  $\mathbf{A}_n(\mathbf{z}) \doteq \sqrt{n}(\mathbf{a}_n(\boldsymbol{\mu}_n + n^{-1/2}\mathbf{z}) - \mathbf{a}_n(\boldsymbol{\mu}_n))$  is a discrete derivative of  $\mathbf{a}_n$  at  $\boldsymbol{\mu}_n$  in the direction  $\mathbf{z}$ . If  $\boldsymbol{\mu}_n$  converges to the fixed point  $\mathbf{f}_k^0$  of the fluid limit mentioned earlier for some  $k \geq 1$ , then  $\sqrt{n}\mathbf{M}_n$  converges to a one-dimensional Brownian motion along the unit vector  $\mathbf{e}_k$  (Lemma 27). However, the standard techniques for proving tightness of  $\mathbf{Z}_n$  from (1.2) do not work, because certain terms of  $\mathbf{A}_n$  can diverge as  $d_n \rightarrow \infty$ . Instead, by carefully studying the asymptotic behavior of  $\mathbf{A}_n$  for suitable regimes of  $\lambda_n \rightarrow 1$  and  $d_n \rightarrow \infty$ , we will directly show that  $\mathbf{Z}_n$  converges to the unique solution of a certain stochastic differential equation. Under suitable initial conditions,  $(Z_{n,k+2}, Z_{n,k+3}, \dots)$  converges to zero in  $l_2$ , and hence let us focus on the limit of the finite dimensional system  $(Z_{n,1}, Z_{n,2}, \dots, Z_{n,k+1}) \Rightarrow (Z_1, Z_2, \dots, Z_{k+1})$ . We find several interesting limiting behaviors:

- (a)  $1 \ll d_n \ll {}^{k+1}\sqrt{n}$  and  $d_n^k(1 - \lambda_n) = \log d_n + o(\sqrt{d_n})$ : Then  $(Z_1, \dots, Z_{k-1}) = 0$ ,  $(Z_k, Z_{k+1})$  is a two dimensional linear diffusion driven by a one-dimensional Brownian motion, and  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_k^0$ . If  $k = 1$ , the limiting behavior of  $\mathbf{A}_n$  is linear. If  $k > 1$  the first  $k$  coordinates of  $\mathbf{A}_n$  diverge, but carefully studying (1.2) shows that this forces the first  $k - 1$  coordinates of  $\mathbf{Z}_n$  to quickly converge to zero (even if they don't start at zero at time 0), while  $\mathbf{A}_n$  behaves linearly for coordinates  $k + 1$  onward.
- (b)  $d_n \sim c\sqrt{n}$  and  $1 - \lambda_n = \frac{\log d_n}{d_n} + \frac{\alpha}{\sqrt{n}}$  for  $\alpha \in (-\infty, \infty)$ : Then  $k = 1$ ,  $\boldsymbol{\mu}_n \rightarrow \mathbf{e}_1$ , and  $(Z_1, Z_2)$  is a two dimensional non-linear diffusion driven by a one-dimensional Brownian motion. In this

case the limit of the first two coordinates of  $\mathbf{A}_n(\mathbf{Z}_n)$  in (1.2) involve the non-Lipschitz term  $\exp(cZ_1)$ , but since the exponential term opposes the growth of the system, a unique pathwise solution  $(Z_1, Z_2)$  for the limiting diffusion equation still exists.

- (c)  $\sqrt{n} \ll d_n \leq n$  and  $1 - \lambda_n = \frac{\log d_n}{d_n} + \frac{\alpha}{\sqrt{n}}$  for  $\alpha \in (0, \infty)$ : Then  $k = 1$ ,  $\boldsymbol{\mu}_n \rightarrow \mathbf{e}_1$ , and  $(Z_1, Z_2)$  is a constrained diffusion on  $(-\infty, \alpha] \times \mathbb{R}$  driven by a one-dimensional Brownian motion. Under suitable conditions  $\alpha = \infty$  is allowed too, in which case there is no reflection. When  $\alpha < \infty$ , the first coordinate of  $\mathbf{A}_n(\mathbf{Z}_n)$  is negligible if  $Z_{n,1} < \alpha$  and strongly negative if  $Z_{n,1} > \alpha$ . Thus by studying the excursion of  $Z_{n,1}$  above  $\alpha$ , we will prove that the Skorokhod map on  $\mathbb{R}$  with a downward reflection at  $\alpha \in \mathbb{R}$  emerges in the limit.

See Theorems 9, 10, 11 in Chapter 7 for the precise statement of the above results. Note that Theorem 9 has the most general form since it allows  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_k^0$  for any  $k \geq 1$ . Here  $k$  may be interpreted as the average time spent by a job in the system since  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_k^0$  shows that all but a vanishingly small fraction of queues have exactly  $k$  customers in the near-equilibrium state. In contrast, like most of the results in the existing literature (see Table 1.1), Theorems 10 and 11 only allow the case of  $k = 1$ .

### 1.2.1 Organization of the second part of this dissertation

In Chapter 5 we provide a literature review for the balls and bins problem and the Supermarket model. Then, as a warm up, in Chapter 6 we study the maximum bin load for the balls and bins problem as  $d = d_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Finally in Chapter 7, we present our main results on limit theorems for the Supermarket model under heavy traffic as  $(d_n, \lambda_n) \rightarrow \infty$  as  $n \rightarrow \infty$ .

## CHAPTER 2

### Literature review for bimodule detection

In this chapter we will review literature from the areas of correlation mining and multi-view data analysis that will be relevant to the bimodule detection problem discussed in Chapter 1.

#### 2.1 Introduction

As discussed in Chapter 1, in numerous scientific disciplines like genomics, neuroscience, ecology and climate science, one frequently has multiple types of measurements obtained from the same underlying samples. The different measurement types, sometimes called data views or modalities, may interact and inform each other, but they typically capture complementary information about the underlying phenomenon. Hence, novel statistical techniques have emerged that jointly analyze data from different measurement modalities to understand the relationship between modalities, distinguish samples based on their joint behavior and achieve other objectives. This collection of techniques is known by various names like data integration, multi-view or multi-modal data analysis, or multi-table methods. Compared to analyzing each measurement modality separately, multi-view data analysis has the potential to improve statistical power and to allow us to answer new types of questions (68, 113).

For multi-view data, we are interested in the exploratory problem of finding bimodules. Recall that the pair  $(A, B)$  is called a bimodule if  $A$  and  $B$  are sets of features from two data types, called Type 1 and Type 2, so that the aggregate cross-correlation between features from  $A$  and  $B$  is statistically significant. This problem falls under the general area of correlation mining, which aims to infer various aspects about the population covariance matrices based on observed samples. In the modern setting of high dimensional features and limited sample size, is typically not feasible to accurately estimate the entire population covariance (or the related correlation or precision) matrix, but recent developments have shown that tasks related to population covariance matrices



such as model selection and screening (59), and various local and global hypothesis tests (27) can still be performed in a statistically principled way.

For the purpose of finding bimodules, we might hope to apply the recent ideas surveyed in (27, 59) to the cross-correlation matrix of Type 1 and Type 2 features to find entries of the matrix that are significant. This is similar to the traditional eQTL analysis method discussed in Chapter 1, albeit asymptotic theoretical analysis of false discoveries could further be used to improve power (28, 58, 57). However, recall that bimodules detect groups of cross-correlated features rather than individual pairs, so the output of pairwise procedures would further need to be analyzed to produce groups (see network based methods below).

## 2.2 Survey of methods for bimodule detection

In this section, we will consider multi-view data analysis and correlation mining methods that can be used to produce bimodules. Many of these methods have been developed or applied in the context of eQTL analysis.

### 2.2.1 Canonical Correlation Analysis

Canonical Correlation Analysis (CCA) is one of the oldest methods for multi-view data analysis, invented by Hotelling in 1936 (60). CCA has been used in various biomedical, environmental, and genomic tasks (see references in (128)). Given two data types, CCA discovers pairs of maximally correlated vectors that are linear combinations of features from each of the data types. These linear combinations, called *canonical variates*, represent latent features that capture aggregate association between the two data types. Multiple canonical variates are discovered by enforcing the constraint that different variates be uncorrelated. An appealing theoretical feature of CCA is that these canonical variates can directly be calculated from the spectral decomposition of  $\Sigma_1^{-1/2}\Sigma_{12}\Sigma_2^{-1/2}$  where  $\Sigma_{12}$  is the cross-covariance matrix of the two data types and  $\Sigma_i$  for  $i = 1, 2$  are the intra-covariance matrices of the two data types (see 90, Chapter 11). Hence, at least when these covariance matrices can be accurately estimated, CCA aggregates signal and accounts for the effects of intra-correlations.

But there are problems with using CCA in the high-dimension setup. The traditional CCA theory relies on accurate estimation of the covariance matrices and their inverses, which may be error-prone when the sample size is not sufficiently large. In fact, if the sample size is less than the number of features, sample covariance matrices will not even be invertible. Secondly, linear combinations of a large number of features may be difficult to interpret. To surpass both these limitations, sparse Canonical Correlation Analysis (sCCA), which restricts attention to canonical variates that are linear combination of only a few features from both the types, was introduced in (101). Importantly for us, each such sparse canonical variate pair can be regarded as a bimodule consisting of the features that appear as linear combinations in the pair.

There have been various further improvements to sCCA (125, 71, 130, 102, 128). Notably, extending the formulation of CCA in terms of an alternative regression scheme, (128) is able to successfully incorporate the intra-covariance matrices in the search for sparse canonical variates. However, in Chapter 4 of this dissertation, we use the sCCA software (129) based on (130), which for computational considerations assumes that the population intra-correlation matrix is identity.

### 2.2.2 Community detection in bipartite networks

Since bimodules are defined in terms of cross-correlations, it is natural to investigate them in the context of the bipartite *cross-correlation* network, which is formed by connecting pairs of features from different data types with a weighted edge, where the weight is equal to the square of the (sample or population) cross-correlation between the features. One might then use community detected methods (10, 9, 72, 39, 105) for weighted and unweighted bipartite networks to find bimodules. However with any such method, care must be taken to avoid discovering bimodules that may be driven by potentially spurious associations.

Along the above lines, the method CONDOR (107) identifies bimodules by applying a fast community detection method (based on bipartite modularity) to the unweighted bipartite graph obtained by thresholding the sample cross-correlations. In (107, 47), the authors use CONDOR to study properties of eQTL-networks, and show that bimodules can be used to place genes and SNPs into a functional context. Their implementation (106) may need the estimated SNP-gene bipartite graph to have a giant connected component, which may require the inclusion of potentially false edges.

The approach that we will take in Chapter 3 will be network based, but it will differ from community-detection based approaches such as CONDOR. While stable population bimodules will be defined in terms of the population cross-correlation network, the sample cross-correlation network will *not* be a sufficient statistic for stable sample bimodules, which will also account for the intra-correlations between features of the same type.

### 2.2.3 Gaussian models, penalized regression, and other approaches

One of the standard techniques to discover conditional relationships between features in moderate- to high-dimensional datasets is to fit a sparse Gaussian graphical model (GGM) to the data by maximizing a penalized version of the Gaussian likelihood (59). For instance, methods such as (132, 134, 36) have been used to uncover conditional dependencies in the eQTL network. However, more relevant to the problem of detecting bimodules, Cheng et al. (34, 32, 33) fit sparse linear GGMs to the eQTL data with a hidden layer that model interactions between sets of genes and sets of SNPs. Their original methodology from (34) is improved in (32, 33) to allow for the discovery of both individual and group effects between SNPs and genes, and to allow for a variable number of bimodules.

Along another direction, many papers (see (119) and references therein) have formulated the eQTL detection problem in terms of a penalized multi-task regression. In this setup, the gene matrix is considered as a collection of response variables (note the term ‘multi-task’) that is jointly regressed against the SNP matrix. In addition to enforcing element-wise sparsity on the matrix of gene-SNP regression coefficients, the penalization term may also be used to enforce various group sparsity and fusion constraints that leverage additional information like the intra-gene and intra-SNP interaction networks to improve detection power. Notably, in the two-graph guided multi-task lasso (31), the gene and SNP networks derived from the respective intra-correlations matrices are used to create a fusion penalty that enforces a network-to-network mapping in the same spirit as bimodules. However, it is not clear if this method is accounting for the intra-correlations in the same way as we intend: i.e. to attenuate significance of cross-correlations due to the presence of intra-correlations.

Finally, one may also search for bimodules by applying a standard clustering method such as k-means to a joint data matrix containing standardized features from the two data types, and

treating any cluster with features from both data types as a bimodule. While appropriate as a “first look”, this approach requires specifying the number of clusters, imposes the constraint that every feature be part of one, and only one, bimodule, and, most importantly, does not distinguish between cross- and intra-correlations.

## CHAPTER 3

### Bimodule Search Procedure

In this chapter, we will introduce the Bimodule Search Procedure that searches for *stable* bimodules based on iterative testing. We will first, in Section 3.2, define the notion of a stable bimodule at the population level, followed by, in Section 3.3.2, its definition at the sample level. The relationship of stable bimodules to bipartite networks and Nash Equilibria stated along the way will motivate the Bimodule Search Procedure described in Section 3.3.3. We now begin with notation and basic assumptions.

#### 3.1 Notation and stochastic setting

Suppose that we have acquired data sets of two different types from a common set of  $n$  samples. Let  $\mathbb{X}$  be an  $n \times p$  matrix containing the data of Type 1, and let  $\mathbb{Y}$  be an  $n \times q$  matrix containing the data of Type 2. The  $i$ th row of  $\mathbb{X}$  and  $\mathbb{Y}$  contain the measurements of Type 1 and Type 2, respectively, on the  $i$ th sample. The columns of  $\mathbb{X}$  and  $\mathbb{Y}$  correspond to the measured features of each type. Denote features of Type 1 by  $S = \{s_1, s_2, \dots, s_p\}$ , and those of Type 2 by  $T = \{t_1, t_2, \dots, t_q\}$ . We assume that the rows of the joint matrix  $[\mathbb{X}, \mathbb{Y}]$  are independent copies of a random (row) vector

$$(\mathbf{X}, \mathbf{Y}) = (X_{s_1}, \dots, X_{s_p}, Y_{t_1}, \dots, Y_{t_q}).$$

For each  $s \in S$  let  $\mathbb{X}_s$  be the column of  $\mathbb{X}$  corresponding to feature  $s$ ; for each  $t \in T$  let  $\mathbb{Y}_t$  be the column of  $\mathbb{Y}$  corresponding to feature  $t$ . For  $s \in S$  and  $t \in T$  let  $\rho(s, t)$  be the population correlation between the random variables  $X_s$  and  $Y_t$ , and let  $r(s, t)$  denote the sample correlation between  $\mathbb{X}_s$  and  $\mathbb{Y}_t$ . For  $A \subseteq S$  and  $B \subseteq T$  we define the aggregate squared (population and sample) correlation

between  $A$  and  $B$  by

$$\rho^2(A, B) \doteq \sum_{s \in A, t \in B} \rho^2(s, t), \text{ and} \quad (3.1)$$

$$r^2(A, B) \doteq \sum_{s \in A, t \in B} r^2(s, t). \quad (3.2)$$

For singleton sets we will omit brackets, writing  $\rho^2(s, B)$  and  $\rho^2(A, t)$  instead of  $\rho^2(\{s\}, B)$  and  $\rho^2(A, \{t\})$ .

### 3.2 Stable population bimodules

Recall that our goal is to identify pairs  $(A, B)$  with  $A \subseteq S$  and  $B \subseteq T$  such that the aggregate cross-correlation between features in  $A$  and  $B$  is large. We begin our analysis of this problem at the population level, where the aggregate cross-correlation between  $A$  and  $B$  can be measured by the quantity  $\rho^2(A, B)$  defined in (3.1). One might rank pairs  $(A, B)$  using a score based on  $\rho^2(A, B)$ , but it is not immediately clear how such a score should be defined. For example,  $\rho^2(A, B)$  itself favors larger feature sets, while the average  $\rho^2(A, B)/|A||B|$  might favor smaller feature sets. More importantly, it is not immediately clear how a score based on  $\rho^2(A, B)$  can be effectively translated to, and evaluated in, the sample setting.

To address these issues, we shift our assessment of pairs  $(A, B)$  from global numerical performance measures to internal stability criteria that are based on the structure of the population cross-correlations. The basic idea is contained in the following definition.

**Definition 1.** A pair  $(A, B)$  of non-empty sets  $A \subseteq S$  and  $B \subseteq T$  is a *stable population bimodule* if

1.  $A = \{s \in S \mid \rho^2(s, B) > 0\}$  and
2.  $B = \{t \in T \mid \rho^2(A, t) > 0\}$ .

In words, the definition says that  $A$  is exactly the set of features in  $S$  that are correlated in aggregate with the features in  $B$ , while  $B$  is exactly the set of features in  $T$  that are correlated in aggregate with the features in  $A$ . It is useful to consider stable bimodules in the context of the population network of cross-correlations.

**Definition 2.** The *population cross-correlation network*  $G_p$  is the weighted bipartite network with vertex set  $S \cup T$ , edge set  $E_p = \{(s, t) \in S \times T \mid \rho(s, t) \neq 0\}$ , and weight function  $w_p : E \rightarrow [-1, 1]$  given by  $w_p(s, t) = \rho(s, t)$ .

The following elementary lemma shows that stable bimodules are closely related to the connected components of  $G_p$ .

**Lemma 1.** A pair  $(A, B)$  of non empty sets with  $A \subseteq S$  and  $B \subseteq T$  is a population bimodule if and only if  $A \cup B$  is a union of non-trivial connected components of  $G_p$ .

*Proof.* For any subsets  $F \subseteq S$  and  $G \subseteq T$  note that  $\rho^2(F, G) > 0$  if and only if  $(F \times G) \cap E_p \neq \emptyset$ . Let  $\text{Nb}(s) \doteq \{t' \mid (s, t') \in E_p\}$  and  $\text{Nb}(t) \doteq \{s' \mid (s', t) \in E_p\}$  and denote the neighborhoods of  $s$  and  $t$  in the graph  $G_p$ . The two conditions in Definition 1 are equivalent to saying  $A = \cup_{t \in B} \text{Nb}(t)$  and  $B = \cup_{s \in A} \text{Nb}(s)$ , respectively. Equivalently, since  $G_p$  is a bipartite graph, the set of nodes  $H = A \cup B$  satisfies the property

$$H = \text{Nb}(H) \tag{3.3}$$

where  $\text{Nb}(C) \doteq \cup_{v \in C} \text{Nb}(v)$  for any subset of vertices  $C \subseteq S \cup T$ .

Using (3.3), let us show that  $H$  is a union of non-trivial connected components. For any  $r \in S \cup T$ , note that the connected component containing  $r$  is defined by  $C(r) \doteq \cup_{i=0}^{\infty} C_i$  where  $C_0 = \{r\}$  and  $C_i = \text{Nb}(C_{i-1})$  for each  $i \geq 1$ . For  $r \in H$ , repeatedly applying (3.3) shows that  $r \in C(r) \subseteq H$  and hence

$$H = \cup_{r \in H} C(r). \tag{3.4}$$

Since (3.3) holds and  $G_p$  is a simple graph, each  $r \in H$  has at least one other neighbor. Hence  $|C(r)| > 1$  for all  $r \in H$ .

Finally, since  $\text{Nb}(C(r)) \subseteq C(r)$  for any  $r$ , if  $H$  satisfies (3.4) then  $\text{Nb}(H) \subseteq H$ . Moreover if all the connected components in (3.4) are non-trivial then  $\text{Nb}(H) \supseteq H$  and (3.3) is satisfied.  $\square$

As the lemma shows, stable population bimodules depend only on the edges of  $G_p$ ; they do not depend on the edge weights, or on correlations between features of the same type. As we will see below, the situation for sample bimodules is substantially different.

### 3.2.1 Bimodules are Nash Equilibria

The notion of stability in Definition 1 has close connections with Nash equilibrium (93) in game theory. To make this precise, fix an  $\epsilon > 0$ , and consider the reward function  $\Phi_\epsilon$  that for any  $A \subseteq S$  and  $B \subseteq T$  takes the value

$$\Phi_\epsilon(A, B) \doteq \sum_{s \in A} \sum_{t \in B} \rho^2(s, t) - \epsilon |A| |B|. \quad (3.5)$$

Consider a two player game in which player 1 chooses a subset  $A \subseteq S$ , player 2 chooses a subset  $B \subseteq T$ , and the payoff to both the players is  $\Phi_\epsilon(A, B)$ . In this setting, a pair of subsets  $(A^*, B^*)$  is called a Nash equilibrium if

$$\max_{A \subseteq S} \Phi_\epsilon(A, B^*) = \Phi_\epsilon(A^*, B^*) = \max_{B \subseteq T} \Phi_\epsilon(A^*, B).$$

The following elementary lemma shows that bimodules are the just Nash equilibria in this game.

**Lemma 2.** *Let  $\delta = \min \{ \rho^2(s, t) \mid s \in S, t \in T, \rho(s, t) \neq 0 \}$  and  $\epsilon_0 = \delta(\max(|S|, |T|))^{-1}$ . If  $\epsilon \in (0, \epsilon_0)$  then the non-empty Nash equilibria of the game with reward function  $\Phi_\epsilon$  coincides with the family of stable population bimodules.*

*Proof.* Suppose  $0 < \epsilon < \epsilon_0$ . Fix  $A \subseteq S$  and observe that for any  $B \subseteq T$

$$\begin{aligned} \Phi_\epsilon(A, B) &= \sum_{t \in B} \sum_{s \in A} (\rho^2(s, t) - \epsilon) \\ &= \sum_{t \in B} (\rho^2(A, t) - \epsilon |A|). \end{aligned} \quad (3.6)$$

Since  $\epsilon |A| < \epsilon_0 |A| \leq \delta$ , for any  $t \in T$ , if  $\rho^2(A, t) > 0$  then  $\rho^2(A, t) - \epsilon |A| > 0$ . Hence the maximum over subsets  $B \subseteq T$  will be uniquely attained at

$$\arg \max_{B \subseteq T} \Phi_\epsilon(A, B) = \{t \in T \mid \rho^2(A, t) - \epsilon |A| > 0\} = \{t \in T \mid \rho^2(A, t) > 0\} \quad (3.7)$$



Similarly, if we fix  $A \subseteq T$ , we can show

$$\arg \max_{A \subseteq S} \Phi_\epsilon(A, B) = \{s \in S \mid \rho^2(s, B) > 0\} \quad (3.8)$$

Hence the pair of non-empty sets  $(A^*, B^*)$  is a Nash equilibrium if and only if

$$A^* = \{s \in S \mid \rho^2(s, B^*) > 0\} \text{ and } B^* = \{t \in T \mid \rho^2(A^*, t) > 0\}.$$

These conditions are the same as those required for  $(A^*, B^*)$  to be a population bimodule (Definition 1). □

The connection between stable bimodules and two player games in the population setting suggests a simple iterative scheme to find stable bimodules when the cross-correlations  $(\rho(s, t))_{s \in S, t \in T}$  are known. Begin by fixing a non-empty subset  $B_0 \subseteq T$  and any  $\epsilon \in (0, \epsilon_0)$ , where  $\epsilon_0$  is chosen as in Lemma 2. Then repeatedly update  $A_{k+1} = \arg \max_A \Phi_\epsilon(A, B_k)$  and  $B_{k+1} = \arg \max_B \Phi_\epsilon(A_{k+1}, B)$  for  $k \geq 0$ . As the value of the objective function strictly increases at every update, the sets  $(A_k, B_k)$  are guaranteed to converge to a Nash equilibrium after finitely many steps. If  $B_0$  is such that  $\Phi_\epsilon(A_1, B_0) > 0$ , then the Nash equilibrium will be non-empty and hence a population bimodule.

It is illustrative to view this iterative update procedure in terms of the cross-correlation graph  $G_p$ . The proof of Lemma 2 shows that the update steps are equivalent to

$$A_{k+1} = \{s \in S \mid \rho^2(s, B_k) > 0\} \text{ and } B_{k+1} = \{t \in T \mid \rho^2(A_{k+1}, t) > 0\}. \quad (3.9)$$

In other words  $A_{k+1}$  is set of neighbors of  $B_k$ , and  $B_{k+1}$  is the set of neighbors of  $A_{k+1}$ . If the iterative update procedure begins from a singleton set  $B_0 = \{t\}$  for some  $t \in T$ , then it corresponds to the breadth first search algorithm for finding the connected component of  $t$  in  $G_p$  (see, e.g., 38). This connection shows that consideration of singleton sets  $B_0 = \{t\}$  for  $t \in T$  finds all the connected components of  $G_p$ , which by Lemma 1, are the minimal stable population bimodules.

### 3.3 Sample stable bimodules and the Bimodule Search Procedure (BSP)

We now extend the notion of a stable bimodule and the iterative search procedure described above to the sample (empirical) setting using ideas and methods from multiple testing. While the empirical setting involves a number of additional complications, the motivation behind stability is essentially the same.

In practice, the population cross-correlations  $\rho(s, t)$  are unknown, and the search for bimodules is based on the observed data matrices  $[\mathbb{X}, \mathbb{Y}]$ . One may simply replace the population correlations with their sample counterparts  $r(s, t)$  in Definition 1, but when working with continuous data  $r(s, t) \neq 0$  (even if  $\rho(s, t) = 0$ ), and in this case the only stable bimodule is the full index set  $(S, T)$ . To address this, we replace the conditions  $\rho^2(s, B) > 0$  by  $r^2(s, B) > \hat{\gamma}(s, B)$ , where the threshold  $\hat{\gamma}(s, B)$  is derived from the application of an FDR-controlling multiple testing procedure to approximate permutation p-values for the statistics  $\{r^2(s, B) : s \in S\}$ . An analogous approach is taken for the conditions  $\rho^2(A, t) > 0$ .

#### 3.3.1 Permutation null distribution and p-values

**Definition 3.** Let  $[\mathbb{X}, \mathbb{Y}]$  be given, and let  $P_1, P_2 \in \{0, 1\}^{n \times n}$  be chosen independently and uniformly from the set of all  $n \times n$  permutation matrices. The *permutation null distribution* of  $[\mathbb{X}, \mathbb{Y}]$  is the distribution of the data matrix

$$[\tilde{\mathbb{X}}, \tilde{\mathbb{Y}}] \doteq [P_1 \mathbb{X}, P_2 \mathbb{Y}].$$

Let  $\mathbf{P}_\pi$  and  $\mathbf{E}_\pi$  denote probability and expectation, respectively, under the permutation null. For  $s \in S$  and  $t \in T$  let  $R(s, t)$  be the (random) sample-correlation of  $\tilde{\mathbb{X}}_s$  and  $\tilde{\mathbb{Y}}_t$ .

The permutation null distribution is obtained by randomly reordering the rows of  $\mathbb{X}$  and independently doing the same for the rows of  $\mathbb{Y}$ . Permutation preserves the sample correlation between features in  $S$ , and between features in  $T$ , but it nullifies the cross-correlations between features in  $S$  and  $T$ . Indeed, as shown in (136),  $\mathbf{E}_\pi[R(s, t)] = 0$  for each  $s \in S$  and  $t \in T$ .

**Definition 4.** For  $A \subseteq S$  and  $B \subseteq T$  define the permutation p-value

$$p(A, B) \doteq \mathbf{P}_\pi(R^2(A, B) \geq r^2(A, B)) \quad (3.10)$$

where  $R^2(A, B) \doteq \sum_{s \in A, t \in B} R^2(s, t)$  and the observed sum of squares  $r^2(A, B)$  is fixed.

The permutation p-value  $p(A, B)$  is the probability under the permutation null distribution that the aggregate cross-correlation between the features in  $A$  and  $B$  exceeds its observed value in the data. Small values of  $p(A, B)$  provide evidence in favor of the hypothesis that  $\rho^2(A, B) > 0$ . As the permutation distribution preserves the correlations between features from  $A$  and between features from  $B$ ,  $p(A, B)$  accounts for the presence of these correlations while assessing the significance of  $r^2(A, B)$ .

### 3.3.2 Stable sample bimodules

Let  $p = (p_1, \dots, p_m)$  be a sequence of p-values and let  $\alpha \in (0, 1)$  be a target false discovery rate. Recall that the Benjamini–Yekutieli (B.Y.) (12) rejection threshold at level  $\alpha$  is defined by

$$\tau_\alpha(p) \doteq \max \left\{ p_{(j)} : \frac{m p_{(j)}}{j} \leq \frac{\alpha}{\sum_{i=1}^m i^{-1}} \right\} \quad (3.11)$$

where  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$  are the ordered values of  $p$ .

**Definition 5.** (Stable Sample Bimodule) Let  $[\mathbb{X}, \mathbb{Y}]$  and  $\alpha \in (0, 1)$  be given. A pair  $(A, B)$  of non-empty sets  $A \subseteq S$  and  $B \subseteq T$  is a *stable sample bimodule* at level  $\alpha$  if

1.  $A = \{s \in S \mid p(s, B) \leq \tau_\alpha(p_{\cdot, B})\}$  and
2.  $B = \{t \in T \mid p(A, t) \leq \tau_\alpha(p_{A, \cdot})\}$

where  $p_{\cdot, B} = \{p(s, B)\}_{s \in S}$  and  $p_{A, \cdot} = \{p(A, t)\}_{t \in T}$ .

Thus  $(A, B)$  is a sample stable bimodule if  $A$  is exactly the set of features in  $S$  that are significantly correlated in aggregate with the features in  $B$ , and at the same time  $B$  is exactly the set of features in  $T$  that are significantly correlated in aggregate with the features in  $A$ . When no ambiguity will arise, we will refer to stable sample bimodules simply as stable bimodules.

Although the definition above parallels that of Definition 1, critical differences emerge in the sample setting. One key difference is the aggregation of small effects. As noted above that the condition  $p(s, B) \leq \tau_\alpha(p_{\cdot, B})$  is equivalent to requiring  $r^2(s, B) \geq \hat{\gamma}(s, B)$  where  $\hat{\gamma}(s, B)$  depending on  $\tau_\alpha(p_{\cdot, B})$ . The latter condition may be satisfied even if the feature  $s$  is not significantly correlated with any individual feature in  $B$ . Similar remarks apply to  $p(A, t)$

Another, more important, difference between the population and sample settings is the role of intra-correlations. A likely side-effect of any empirical search for pairs  $(A, B)$  with high cross-correlations is that the intra-correlations of the features in  $A$  and  $B$  will also be large, often significantly larger than the intra-correlations of a randomly selected set of features with the same cardinality. Failure to account for inflated intra-correlations can lead to anti-conservative (optimistic) assessments of significance, false discoveries, and oversized feature sets. As noted above, the permutation distribution leaves intra-correlations unchanged, while ensuring that cross-correlations are equal to zero. In this way the permutation p-values  $p(s, B)$  and  $p(A, t)$  directly account for the intra-correlations among features in  $A$  and  $B$ .

### 3.3.3 The Bimodule Search Procedure (BSP)

We adapt the iterative search procedure for population bimodules described at the end of Section 3.2 using the p-value based characterization of sample bimodules in Definition 5. The result is an iterative, testing-based search procedure for stable bimodules. Iterative-testing based procedures have been used in single data-type settings for community detection (98), differential correlation mining (18), and association mining (17). The definition and approximation of the permutation based p-values used here differs substantially from this existing work.

An overview of BSP is given in Algorithm 1. If BSP terminates at a non-empty fixed point, then its output is a stable bimodule at level  $\alpha$ . Unlike its population counterpart, BSP is not guaranteed to terminate in a finite number of steps: as the procedure operates in a deterministic manner, and the number of feature set pairs is finite, BSP will terminate at a (possibly empty) fixed point or enter a limiting cycle. To limit computation time, the loop at Line 2 is stopped after 20 iterations. In our simulations and real-data analyses (described in Chapter 4) the 20 iteration limit was enforced in only a handful of cases. Further details on how BSP deals with cycles and limits large sets can be found in Section 3.4.1.

**Input :** Data matrices  $\mathbb{X}$  and  $\mathbb{Y}$  and parameter  $\alpha \in (0, 1)$ .

**Result:** A stable bimodule  $(A, B)$  at level  $\alpha$ , if found.

```

1 Initialize  $A' = \{s\} \subseteq S$  and  $B' = \emptyset$ ;
2 do
3    $(A, B) \leftarrow (A', B')$ ;
4   Compute  $p(A, t)$  for each  $t \in T$  and let  $p_T \leftarrow (p(A, t))_{t \in T}$ ;
5    $B' \leftarrow \{t \in T \mid p(A, t) \leq \tau_\alpha(p_T)\}$ ; // The indices rejected by B.Y.
6   Compute  $p(s, B')$  for each  $s \in S$  and let  $p_S \leftarrow (p(s, B'))_{s \in S}$ ;
7    $A' \leftarrow \{s \in S \mid p(s, B') \leq \tau_\alpha(p_S)\}$ ; // The indices rejected by B.Y.
8 while  $(A', B') \neq (A, B)$ ;
9 if  $|A||B| > 0$  then
10 | return  $(A, B)$ ;
11 end

```

**Algorithm 1:** Bimodule Search Procedure (BSP)

In practice, BSP is initialized with each singleton pair  $(s, \emptyset)$  for  $s \in S$ , and each singleton pair  $(\emptyset, t)$  for  $t \in T$ . While this initialization guarantees the recovery of all minimal stable bimodules in the population setting, no such guarantees are available in the sample setting. Nevertheless, we have found this initialization strategy to be effective in practice. When either of the sets  $S$  or  $T$  is large, we use additional strategies to speed up computation, like randomly selecting a smaller subset of features for initialization (Section 3.4.2).

The constant  $\alpha \in (0, 1)$  is the only free parameter of BSP. We will refer to  $\alpha$  as the false discovery parameter. While  $\alpha$  controls the false discovery rate at each step of the search procedure, this does not guarantee control on the false associations (i.e.  $(s, t)$  such that  $\rho(s, t) = 0$ ) within the stable bimodules. In general, BSP will find fewer and smaller bimodules when  $\alpha$  is small, and find more numerous and larger bimodules when  $\alpha$  is large. In practice, we employ a permutation based procedure to select  $\alpha$  from a fixed grid of values based on the notion of *edge-error*. See Section 3.4.3 for details.

Simulations and theoretical calculations suggest that singleton bimodules  $(\{s\}, \{t\})$  at a given level  $\alpha \in (0, 1)$  can occur even in completely random data if  $|S|$  and  $|T|$  are large enough. To minimize the detection of spurious singleton bimodules, we discard bimodules  $(A, B)$  with  $p(A, B) > \frac{\alpha}{|S||T|}$ , where the threshold is the Bonferroni correction at level  $\alpha$  for singleton bimodules. Alternatively, one can simply discard singleton bimodules with p-values exceeding the Bonferroni threshold.

The BSP search procedure often finds the same bimodule starting from multiple initializations, and in some cases there are numerous bimodules having substantial overlap. In the latter case, we assess the *effective* number of distinct bimodules and select an equal number of representative bimodules for subsequent analysis. Details can be found in Section 3.4.4.

### 3.3.3.1 Approximation of p-values

Recall that BSP is not based on an underlying generative or distributional model. The method relies on the assumption that the samples are independent and identically distributed and on Definition 3, the permutation based p-values  $p(s, B)$  and  $p(A, t)$ . A total of  $|S| + |T|$  p-values are calculated in each iteration of the loop at Line 2 in Algorithm 1. Accounting for multiple initializations, several billion p-value calculations are required for typical genomic data sets. Moreover, the resolution of these p-values must be small enough to account for multiple-testing correction.

When  $|S|$  or  $|T|$  is large, calculating the p-values  $p(s, B)$  and  $p(A, t)$  using a standard Monte Carlo permutation scheme is not feasible. As an alternative, we make use of ideas from (136) and (137) to approximate the permutation p-values  $p(A, t)$  and  $p(s, B)$  using the tails of a location-shifted Gamma distribution having same first three moments as the sampling distribution of  $R^2(A, t)$  under the permutation null.

Although the first three moments of  $R^2(A, t)$  can be computed exactly (136), to further speed computation we use instead the eigenvalue conditional moments of  $R^2(A, t)$  (see (137), also called the continuous permutation scheme in Section 3.5 below), which depend only on the eigenvalues of the intra-correlation matrix of the features in  $A$ , and not on  $t$ . The analytical formula for the eigenvalue conditional moments is based on a normality assumption for the data generating distribution, but one may show that the weaker assumption of spherical symmetry is sufficient. In practice, the additional assumptions used in the moment approximation do not appear to limit the applicability of BSP. Accuracy of the p-value approximations is briefly discussed in Section 3.4.6.

### 3.3.4 Network interpretation of sample stable bimodules

As discussed in Section 3.2, stable population bimodules can be studied in terms of the correlation network  $G_p$ , and it is useful to study sample bimodules in a similar manner. To this end, we

define the *sample cross-correlation network*  $G_s$  to be the weighted bipartite network with vertex set  $S \cup T$ , full edge set  $E = S \times T$ , and weight function  $w(s, t) = r(s, t)$ .

**Definition 6.** For each  $\tau > 0$  define the (unweighted) network  $G_s^\tau = (S \cup T, E(\tau))$  where  $E(\tau) = \{(s, t) \in S \times T \mid |w(s, t)| \geq \tau\}$ . For each feature set pair  $(A, B)$  we define the *connectivity-threshold*

$$\tau^*(A, B) = \max\{\tau \in [0, 1] : A \cup B \text{ is connected in } G_s^\tau\}. \quad (3.12)$$

It follows from Lemma 1 that minimal population bimodules (those obtained when starting the iterative search from singletons) correspond to connected components of  $G_p$ . Accordingly, we define the *essential-edges* of a bimodule  $(A, B)$  to be those that are present at the connectivity threshold

$$\text{essential-edges}(A, B) = (A \times B) \cap E(\tau^*(A, B)). \quad (3.13)$$

Note that  $E(\tau) \cap (A \times B)$  is a set-estimate of the edges  $(s, t) \in A \times B$  with  $\rho(s, t) \neq 0$ , and that the choice of  $\tau > 0$  affects the fraction of false discoveries in this estimate. The value  $\tau = \tau^*(A, B)$  is the most conservative threshold subject to the constraint that  $A \cup B$  is connected in  $G_s^\tau$ , and the essential edges are those of the resulting graph. Assuming that the bimodule  $(A, B)$  is connected in the population network, we expect the essential-edges to be a conservative estimate of the true edges in the population network.

### 3.4 BSP implementation details

In this section we dive into further implementation details of BSP.

#### 3.4.1 Dealing with cycles and large sets

In practice, we do not want the sizes of the sets  $(A_k, B_k)$  in the iteration to grow too large as this slows computation, and large bimodules are difficult to interpret. Therefore the search procedure is terminated when the geometric size of  $(A_k, B_k)$  exceeds 5000. In some cases, the sequence of iterates  $(A_k, B_k)$  for  $k \in \{1, \dots, k_{max}\}$  will form a cycle of length greater than 1, and will therefore fail to reach a fixed point. To search for a nearby fixed point instead, when we encounter the cycle

$(A_k, B_k) = (A_l, B_l)$  for some  $l < k - 1$ , we set  $(A_{l+1}, B_{l+1})$  to  $(A_k \cap A_{k-1}, B_k \cap B_{k-1})$  and continue the iteration.

### 3.4.2 Initialization heuristics for large datasets

When  $S$  is large, we initialize BSP from all the features in  $T$ , but initialize only from a subset of randomly chosen features in  $S$ . We also note that identical resultant bimodules are repeatedly discovered by BSP starting from different initializations, often from features within said bimodule. This problem is particularly prominent for large bimodules which may be rediscovered by upto several thousand initializations. Hence, to avoid some of this redundant computation, we may skip initializing BSP from features in the bimodules that have already been discovered.

### 3.4.3 Choice of $\alpha$ using half-permutation based edge-error estimates

To select the false discovery parameter  $\alpha$  for BSP, we estimate the *edge-error* for each value of  $\alpha$  from a pre-specified grid. The edge-error is an edge-based false discovery notion for bimodules, defined as the average fraction of erroneous essential-edges (defined in Section 3.3.4) among bimodules. Since we do not know the ground truth, we estimate the edge-error for BSP by running it on instances of the *half-permuted* dataset in which the sample labels for half of the features from each data type have been permuted. Further details are given below.

#### 3.4.3.1 Half-permutation

Comparing results between the original and permuted data (Definition 3) allow us to assess the false discoveries from BSP when there are no true associations between features from  $S$  and  $T$ . However, we often expect associations between at least some variables from  $S$  and  $T$  (in fact, these are the ones that we want to find). To create a null distribution where some of pairs of features from  $S$  and  $T$  are correlated and some are not, we use the following half-permutation scheme. Let  $(\mathbb{X}, \mathbb{Y})$  denote the original data, where  $\mathbb{X}$  and  $\mathbb{Y}$  are measurements matrices for the two data types. We generate a *half-permuted* dataset  $(\tilde{\mathbb{X}}, \tilde{\mathbb{Y}})$  as follows:

1. Randomly select half the features,  $\hat{S} \subseteq S$  and  $\hat{T} \subseteq T$ , from each data type.



2. Randomly permute the rows of the submatrix of  $\mathbb{X}$  that corresponding to the columns  $\hat{S}$ , and call the resulting matrix  $\tilde{\mathbb{X}}$ . In other words, the submatrix corresponding to the features  $S \setminus \hat{S}$  is the same in  $\mathbb{X}$  and  $\tilde{\mathbb{X}}$ , while the sample labels of the submatrix of  $\tilde{\mathbb{X}}$  corresponding to features in  $\hat{S}$  have been permuted  $\mathbb{X}_{\hat{S}} = \mathbf{P}_1 \mathbb{X}_{\hat{S}}$  by a random permutation matrix  $\mathbf{P}_1$ .
3. Similarly, permute the rows of matrix  $\mathbb{Y}$  corresponding to the features  $\hat{T}$  using another independent permutation matrix  $\mathbf{P}_2$ . Call the resulting matrix  $\tilde{\mathbb{Y}}$ .

Note that, together, the “half-permutation” steps 2 and 3 temper the cross-correlation between pairs of features in  $\hat{S} \times T \cup S \times \hat{T}$ . Let  $\mathcal{B} = \{(A_1, B_1), (A_2, B_2) \dots (A_K, B_K)\}$  be the collection of bimodules in the half-permuted data  $(\tilde{\mathbb{X}}, \tilde{\mathbb{Y}})$ . We will assume that the collection  $\mathcal{B}$  is already filtered for overlaps (Section 3.4.4 below). We define the edge-error estimate for the collection  $\mathcal{B}$  as

$$\widehat{\text{edge-error}}(\mathcal{B}) = \frac{1}{|\mathcal{B}|} \sum_{(A,B) \in \mathcal{B}} \frac{\left| \text{essential-edges}(A, B) \cap \left( \hat{S} \times T \cup S \times \hat{T} \right) \right|}{|\text{essential-edges}(A, B)|}. \quad (3.14)$$

In practice, we generate half-permuted datasets and use the edge-error estimate (3.14) to choose  $\alpha$  as follows. First, we generate a pre-specified number  $N$  of instances of the half-permuted dataset. If the covariates are present, we correct for them after the half-permutation step. Next, for each  $\alpha$  among a range of values, e.g.  $\{0.01, 0.02, \dots, 0.05\}$ , we run BSP with false discovery parameter  $\alpha$  over the each of the half-permuted datasets and calculate the average edge-error (3.14) of the resulting collection of bimodules for that value of  $\alpha$ , averaged over all the  $N$  half-permuted instances. We can then choose an  $\alpha$  from the grid that has average edge-error smaller a pre-specified value like 0.05. Generally smaller values of  $\alpha$  tend to have smaller edge error, so we choose the largest value of  $\alpha$  from the grid that has acceptable edge error. However, we may chose a smaller value of  $\alpha$  if the bimodules are too large.

A caveat with the above procedure to select  $\alpha$  is that the edge-error estimates may be quite variable even when averaged over a large number  $N$  of half-permuted datasets. One explanation for this variability is that different  $\hat{S}$  and  $\hat{T}$  are chosen for each instance of the half-permutation. Nevertheless, if we observe variability we choose a more conservative value of  $\alpha$ . As seen in Section 4.1.3 in Chapter 4, even without access to the ground truth, we were able to keep the true edge-error under 0.05 by using the above strategy to select  $\alpha$ .

### 3.4.4 Filtering overlapping bimodules

Running BSP starting from different initializations may lead to a collection  $\mathcal{B} = \{(A_1^*, B_1^*), (A_2^*, B_2^*), \dots (A_J^*, B_J^*)\}$  of  $J$ , potentially repeating and overlapping, bimodules. A typical bimodule might occur multiple times in  $\mathcal{B}$ , and distinct bimodules in  $\mathcal{B}$  might have substantial overlap. To reduce duplication, we count the *effective number* of disjoint bimodules in the collection  $\mathcal{B}$  and propose a way to select a subset of those many bimodules from  $\mathcal{B}$  that are most disjoint. Details follow.

We use the following definition  $N_e(\mathcal{B})$  of *effective number* of disjoint bimodules among the collection  $\mathcal{B}$ , adapted from (115):

$$N_e(\mathcal{B}) \doteq \sum_{(A,B) \in \mathcal{B}} \frac{1}{|A||B|} \sum_{s \in A, t \in B} \frac{1}{C_{\mathcal{B}}(s, t)}, \quad (3.15)$$

where

$$C_{\mathcal{B}}(s, t) = \sum_{(A,B) \in \mathcal{B}} \mathbb{I}_{\{s \in A, t \in B\}} \quad (3.16)$$

is the number of bimodules that contain the pair  $(s, t)$ . As noted in (115), the measure  $N_e$  has the property that if there were  $r$  distinct bimodules in  $\mathcal{B}$ , all disjoint when considered as collection of feature pairs, then  $N_e(\mathcal{B}) = r$ .

When there is substantial overlap between the bimodules in  $\mathcal{B}$ ,  $N_e(\mathcal{B})$  is typically smaller than  $|\mathcal{B}|$ . We can then choose  $N = \lceil N_e(\mathcal{B}) \rceil$  most distinct representative bimodules among  $\mathcal{B}$  using the following two steps:

1. **Cluster the collection of bimodules  $\mathcal{B}$  into  $N$  groups.** We use hierarchical agglomerative clustering using the average linkage method (see 56) based on the distance metric given by

$$d_J((A_1, B_1), (A_2, B_2)) = 1 - \frac{|A_1 \times B_1 \cap A_2 \times B_2|}{|A_1 \times B_1 \cup A_2 \times B_2|} \quad (3.17)$$

When we consider the bimodules as a collection of feature pairs, the metric  $d_J$  is simply the Jaccard distance between the bimodules.

2. **Select one representative bimodule from each cluster.** Let  $\mathcal{C} \subseteq \mathcal{B}$  be one of the  $N$  bimodule clusters obtained in the previous step. We select a bimodule representative

$(A, B) \in \mathcal{C}$  from this cluster that maximizes the following *importance score*:

$$\text{isc}(A, B) = \sum_{s \in A, t \in B} \sum_{(A', B') \in \mathcal{C}} \mathbb{I}_{\{s \in A', t \in B'\}}. \quad (3.18)$$

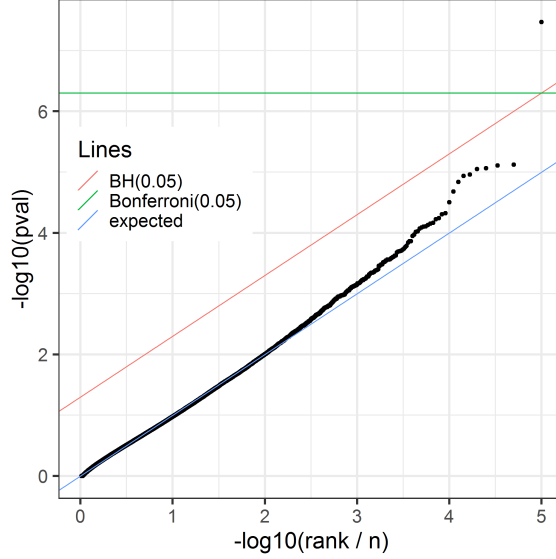
The score  $\text{isc}(A, B)$  aggregates the importance of each feature pair within  $(A, B)$ , measured using the number of (potentially repeating) bimodules from  $\mathcal{C}$  the feature-pair was a part of. Bimodules that are larger and have more overlapping pairs will be preferred as representatives under this scheme.

### 3.4.5 Covariate correction

In some cases the data matrices  $[\mathbb{X}, \mathbb{Y}] \in \mathbb{R}^{n \times (p+q)}$  are accompanied by one or more covariates like sex, platform details and PEER factors (117) (in the case of eQTL data) that must be accounted for by removing their effects before discovering bimodules. Suppose we are given  $m$  such linearly independent covariates  $v_1, \dots, v_m \in \mathbb{R}^n$ . Here we describe how to modify BSP to remove their effects. First, we residualize each column of the original data  $[\mathbb{X}, \mathbb{Y}]$  by setting up a linear model with explanatory variables  $v_1, \dots, v_m$ . Denote the resulting matrix by  $[\mathbb{X}', \mathbb{Y}] \in \mathbb{R}^{n \times (p+q)}$  that has columns which are projections of those of  $[\mathbb{X}, \mathbb{Y}]$  onto the subspace orthogonal to  $v_1, \dots, v_m$ . We would like to now run BSP on  $[\mathbb{X}', \mathbb{Y}]$ , however since the columns of  $[\mathbb{X}', \mathbb{Y}]$  lie on an  $n' = n - m'$  dimensional subspace, the independence assumption of p-value calculations in Section 3.3.3.1 would fail. However, as argued in (136), it is enough to replace the sample size  $n$  with the effective sample size of  $n'$  in the p-value calculations.

### 3.4.6 Checking our p-value approximation

To check the uniformity of our approximate p-value under the permutation null, we chose a bimodule  $(A, B)$  found in Section 4.2 and a  $t \in B$ . Then we randomly permuted the labels of gene  $t$  ( $10^5$  times), computing our p-value approximation  $\hat{p}(A, t)$  in each case. Hence we are assessing the uniformity of  $\hat{p}(A, t)$  under the permutation null distribution. The result in Figure 3.1 shows that the computed p-values are almost uniform but extremely small p-values show anti-conservative behavior. A potential reason for this anti-conservative behavior is that the tails of test statistic under the permutation distribution may be heavier compared to the tails of the location-shifted



**Figure 3.1:** Assessing the accuracy of our p-value estimate  $\hat{p}(A, t)$ : We used the eQTL data from Section 4.2 and chose a bimodule with 24 SNPs (used as  $A$ ) and selected  $t$  to be a gene from the same bimodule. We then performed  $10^5$  random permutation of the sample labels for the gene  $t$  and repeatedly estimated  $\hat{p}(A, t)$  for each permutation after removing the effects of covariates (3.4.5).

Gamma distribution that we use to approximate it, since the permutation distribution is discrete distribution which explicitly depends on the exact entries of the data matrices.

### 3.4.7 Improving computation time

The lines 4 and 6 in Algorithm 1 are usually the most computationally intensive. For instance, the typical steps required to execute lines 6 and 7 would be:

1. Calculate  $r^2(s, B')$  for each  $s \in S$ .
2. Estimate the (null) distribution of  $R^2(s, B')$  for each  $s \in S$  using Definition 3.
3. Obtain the p-value  $p(s) \doteq \mathbf{P}_\pi(R^2(s, B') \geq r^2(s, B'))$  for each  $s \in S$ .
4. Sort the vector  $\mathbf{p} = (p(s))_{s \in S}$  to find the Benjamini-Yekutieli threshold  $\tau_\alpha(\mathbf{p})$ .
5. Return the set  $A' = \{s \in S \mid p(s) \leq \tau_\alpha(\mathbf{p})\}$ .

When  $|S|$  is very large, all the above steps can cause a substantial slowdown. We now describe how these steps could be sped up. Step 1 requires calculating correlations. This involves matrix multiplication that can be done in parallel using BLAS libraries. Since it may be infeasible to store

all possible  $|S| \times |T|$  cross-correlations after they are computed, we store only a limited number of correlations in a Least Recently Used (LRU) cache (95), which the algorithm can utilize when calculating the same set of correlations in succession.

As stated in Section 3.3.3.1, to calculate the distribution of  $R^2(s, B')$  in Step 2, we fit a shifted gamma family to the distribution of  $R^2(s, B')$  based on the first three moments. For further improvement in speed, these moments are calculated using the analytical approximation in (136, 137) for Gaussian data and the continuous permutation scheme described in Section 3.5 below. In more detail, we assume that  $R^2(s, B') \stackrel{d}{\approx} W$ , with

$$W \sim a + c\chi_b^2 \quad (3.19)$$

where  $a, b, c$  are constants chosen so that

$$\begin{aligned} \mathbf{E}W &= \frac{|B'|}{n-1} \\ \text{Var}(W) &= \frac{2}{n^2-1} \left( \sum_i \lambda_i^2 - \frac{|B'|^2}{n-1} \right) \\ \mathbf{E}(W^3) &= \frac{1}{(n^2-1)(n+3)} \left( |B'|^3 + 6|B'| \sum_i \lambda_i^2 + 8 \sum_i \lambda_i^3 \right) \end{aligned}$$

where  $\{\lambda_i\}_{i=1}^{|B'|}$  are the eigenvalues of the intra-correlation matrix of the features in  $B'$ . Note the interesting property that the distribution of  $W$  does not depend on the choice  $s \in S$ . Hence in Step 2, we only need to compute a single distribution. Further using this fact, we now show that Steps 3 – 5 can be performed quickly by computing only a small number of quantiles of  $W$ .

### 3.4.8 Fast Benjamini Hochberg procedure based on quantiles

Suppose  $\mathbf{t} = (t_1, \dots, t_m) \in \mathbb{R}^m$  is an observation of  $m$  (assumed large) test statistics, all of which are hypothesized to have a common null distribution with CDF  $F$ . Hence we assume that the vector of p-values corresponding to these observations is given by  $\mathbf{p} = F(\mathbf{t}) \doteq (F(t_1), \dots, F(t_m)) \in [0, 1]^m$ . Let  $p_{(1)} \leq \dots \leq p_{(m)}$  be the sorted coordinates  $\mathbf{p}$ . The Benjamini Hochberg (11) procedure at level

```

1 Function Fast-BH( $\mathbf{t}, F, \alpha$ ):
    input :  $\mathbf{t} = (t_1, \dots, t_m) \in \mathbb{R}^m$  are observation of  $m$  test statistics,
            $F$  is the common hypothesized null-CDF, and
            $\alpha \in (0, 1)$  is the false discovery parameter.
    output:  $\tau'_\alpha(F(\mathbf{t}))$  – the Benjamini Hochberg threshold for  $\mathbf{p} = F(\mathbf{t})$  at level  $\alpha$ .
2    $i_{new} \leftarrow m, i_{old} \leftarrow \infty$ ;
3   while  $i_{new} < i_{old}$  do
4        $i_{old} \leftarrow i_{new}$ ;
           // The following can be calculated in-place using the data structure
           // used in the quick-sort algorithm.
5        $i_{new} \leftarrow \left| \left\{ j \mid t_j \geq F^{-1}\left(\frac{\alpha i_{old}}{m}\right) \right\} \right|$ ;
           // By induction,  $i_{new} \leq i_{old}$  continues to holds here.
6   end
7   return  $\frac{\alpha i_{new}}{m}$  ;
8 end

```

**Algorithm 2:** Fast Benjamini-Hochberg procedure using quantiles

$\alpha$  defines the rejection threshold

$$\tau'_\alpha(\mathbf{p}) = \sup \left\{ \frac{\alpha i}{m} \mid p_{(i)} \leq \frac{\alpha i}{m}, 1 \leq i \leq m \right\} \vee 0. \quad (3.20)$$

In Algorithm 2, we describe a way to compute  $\tau'_\alpha(\mathbf{p})$  by performing a potentially small number of applications of  $F^{-1}$ . This can be useful when evaluation  $F$  or its inverse  $m$  times is expensive.

### 3.5 The continuous permutation scheme

In the remainder of this chapter, we will discuss a continuous version of the permutation scheme used in Definition 3. This continuous permutation scheme may be of some theoretical interest since it has better analytical properties than the usual permutation scheme, particularly those related to the distribution of correlation matrices. In Sections 3.3.3.1 and 3.4.7, we have used the continuous permutation scheme to approximate the moments of  $R^2(A, B)$  (see Definition 3) with those in Theorem 2 below.

#### 3.5.1 Motivation and notation

Let  $\mathbb{X}$  and  $\mathbb{Y}$  be two data matrices with the same number  $n$  of rows, corresponding to common samples. To search for features (the columns) of  $\mathbb{X}$  and  $\mathbb{Y}$  that are significantly correlated, it is

common to perform hypothesis tests against the *permutation null distribution*  $(\mathbf{P}\mathbb{X}, \mathbb{Y})$ , where  $\mathbf{P}$  is an  $n \times n$  random permutation matrix. The permutation distribution preserves the *intra-correlations* within features from the same matrix, but abates the *cross-correlation* between features from  $\mathbb{X}$  and features from  $\mathbb{Y}$ . In addition to its simplicity, an appeal of the permutation distribution is that it does not make Normality or other parametric assumptions on the data generating distribution. However permutation p-values for even simple statistics, like the sample correlation between a pair of features from  $\mathbb{X}$  and  $\mathbb{Y}$ , do not have closed forms (but their moments may (136)) and often need to be calculated through slower Monte Carlo methods.

One explanation for the mathematical complexity of the permutation distribution is that it is a discrete distribution that depends closely on the exact entries of  $\mathbb{X}$  and  $\mathbb{Y}$ . Here we propose a continuous extension to the permutation distribution that has good mathematical properties. Consider the distribution  $(\mathbf{Q}\mathbb{X}, \mathbb{Y})$  where  $\mathbf{Q}$  is a random  $n \times n$  orthogonal matrix that keeps the  $n$  dimensional vector  $e = (1, 1, \dots, 1)^t$  fixed, generated using the uniform (Haar) measure on the group of such matrices. We will call this the *continuous permutation* scheme because it results in a continuous distribution for the entries of the  $\mathbf{Q}\mathbb{X}$  matrix, unlike under the regular permutation scheme. The columns of the matrix  $\mathbf{Q}$  can be generated by applying the Gram-Schmidt orthogonalization to the sequence of vectors  $\{e, Z_2, \dots, Z_n\}$  where each  $\{Z_i\}_{i=2}^n$  are columns of a  $n \times (n-1)$  matrix with i.i.d.  $\mathcal{N}(0, 1)$  entries.

Like the permutation distribution, the continuous permutation scheme also abates *cross-correlations* in the data while preserving the *intra-correlations*. This can be understood using the following geometric perspective: consider the columns (features) of  $\mathbb{X}$  and  $\mathbb{Y}$  as vectors in the  $n$ -dimensional space. Then the sample (Pearson) correlation between any two features is the inner product between their corresponding columns after projection onto the subspace  $e^\perp \doteq \{v \in \mathbb{R}^n \mid \langle v, e \rangle = 0\}$ . Since  $\mathbf{Q}$  is an orthogonal matrix that keeps the vector  $e$  fixed,  $\mathbf{Q}$  induces a random orthogonal transformation on the space  $e^\perp$ . Hence  $\mathbf{Q}$  applied to features of  $\mathbb{X}$  preserves the intra-correlations, but the transformed features of  $\mathbb{X}$  and features of  $\mathbb{Y}$  no longer tend to be aligned in the projected space, and hence their cross-correlation is reduced. Mathematically, one can show that expected sample correlation between a features from  $\mathbf{Q}\mathbb{X}$  and a feature from  $\mathbb{Y}$  is 0.

In the usual permutation scheme  $\mathbf{Q}$  could only be a permutation matrix, but as seen above, we can use any orthogonal matrix as long as it keeps the vector  $e$  fixed. In this way, the continuous permutation scheme naturally extends the usual permutation scheme since it may be shown that the group of orthogonal matrices that fixes  $e$  is the smallest compact Lie group of positive dimension that contains the subgroup of permutation matrices (61).

As we show below, the distribution of statistics like the sum of squared cross-correlations have analytical expressions under the continuous permutation scheme. These calculations, motivated by Normality assumption on the data, have partly been described in (136, 137) to approximate the permutation distribution of the corresponding statistics. Here we emphasize that Normality assumption on the data is not necessary as long as we work with continuously permuted data. We use the following notation:  $\mathcal{O}_k$  will denote the space of  $k \times k$  orthogonal matrices for any  $k \geq 1$ , and  $e = \frac{1}{\sqrt{n}}(1, \dots, 1) \in \mathbb{R}^n$  throughout this chapter.

### 3.5.2 Distributional results

Let us assume that  $X \in \mathbb{R}^{n \times p}$  and  $Y \in \mathbb{R}^{n \times q}$  are given data matrices. Let  $\mathcal{G}_n = \{U \in \mathcal{O}_n \mid \mathbf{Q}e = e\}$  be the collection of orthogonal matrices that fixes the vector  $e$ . Note that  $\mathcal{G}_n \cong \mathcal{O}_{n-1}$  since elements in  $\mathcal{G}_n$  are determined by their action on the  $n - 1$  dimensional space  $e^\perp = \{v \in \mathbb{R}^n \mid \langle v, e \rangle = 0\}$ . We will prove the following two Theorem for the distribution of sums of squared correlations under the continuous permutation scheme.

**Theorem 1.** *Let  $[\tilde{X}, \tilde{Y}] \doteq [\mathbf{Q}'X, Y] \in \mathbb{R}^{n, p+q}$ , where  $\mathbf{Q}' \in \mathcal{G}_n$  is randomly chosen according to the Haar measure. Let  $\tilde{R}$  denote the  $p \times q$  sample cross-correlation matrix between features (columns) from  $\tilde{X}$  and  $\tilde{Y}$ . Then*

$$T \doteq \sum_{i=1}^p \sum_{j=1}^q \tilde{R}_{i,j}^2 \stackrel{d}{=} \lambda_X^t(\mathbf{Q} \odot \mathbf{Q}) \lambda_Y \quad (3.21)$$

where, for  $A \in \{X, Y\}$ ,  $\lambda_A \in \mathbb{R}^{n-1}$  is the vector of non-zero eigenvalues of the sample correlation matrix of the features of  $A$  (followed by zeros as necessary),  $\mathbf{Q}$  is a random matrix distributed according to the Haar measure on  $\mathcal{O}_{n-1}$ , and  $\mathbf{Q} \odot \mathbf{Q}$  is the matrix with entries given by the square of the corresponding entries of  $\mathbf{Q}$ .

The above theorem says that the distribution of  $T$  – the sum of squares of the entries of the cross-correlation matrix of the continuously permuted data – only depends on the original data



matrices through the eigenvalues of the correlation matrices of  $X$  and  $Y$ . It is worth noting that this is not the case under the usual permutation distribution, which depends more explicitly on the entries of  $X$  and  $Y$ .

The previous theorem provides a distributional equality for the test statistic  $T$ . It shows the distribution of  $T$  is a weighted sum of squares of the entries of the matrix  $\mathbf{Q}$ , but this distribution may not have a standard form. However, the distributional equality can be used to compute up to three moments of test statistic  $T$  as described in the next theorem. These moments can then be used to fit a distribution form the the shifted gamma family to  $T$ .

To motivate the next theorem, let us use (3.21) to calculate the first moment of  $T$ . Denoting the elements of  $\mathbf{Q}$  by  $Q_{ij}$  we see

$$\mathbf{E}(T) = \sum_{i=1}^{n-1} \sum_{j=1}^{n-1} \lambda_{X,i} \lambda_{Y,j} \mathbf{E}(Q_{ij}^2) = \left( \sum_i \lambda_{X,i} \right) \left( \sum_j \lambda_{Y,j} \right) \mathbf{E}(Q_{11}^2) = \frac{pq}{n-1}.$$

where we have used that  $\mathbf{E}(Q_{i,j}^2) = \mathbf{E}(Q_{1,1}^2) = \frac{1}{n-1}$  and  $\sum_i \lambda_{X,i} = p$  and  $\sum_j \lambda_{Y,j} = q$ . Computing higher moments  $E(T^2)$  and  $E(T^3)$  requires knowledge of joint moments of  $\mathbf{Q}$ . These calculations are done in the following theorem.

**Theorem 2.** *Let  $T$  be as defined in Theorem 1 equation (3.21). Let  $s_{A,r} \doteq \sum_{i=1}^{n-1} \lambda_{A,i}^r$  for  $r = 1, 2, 3$  and  $A \in \{X, Y\}$ , and  $N = n - 1$ . Then*

$$\begin{aligned} \mathbf{E}(T^2) &= \alpha s_{X,1}^2 + (\beta - \alpha) s_{X,2} \\ \alpha &= (\alpha_1 - \alpha_2) s_{Y,2} + \alpha_2 s_{Y,1}^2 \\ \beta &= (\beta_1 - \beta_2) s_{Y,2} + \beta_2 s_{Y,1}^2 \end{aligned} \tag{3.22}$$

where  $\alpha_1 = \beta_2 = \frac{1}{N(N+2)}$ ,  $\alpha_2 = \frac{N+1}{(N-1)N(N+2)}$  and  $\beta_1 = \frac{3}{N(N+2)}$ . And

$$\begin{aligned} \mathbf{E}(T^3) &= \alpha s_{X,1}^3 + \{c - a - 3(b - a)\} s_{X,3} + 3(b - a) s_{X,2} s_{X,1} \\ a &= a_3 s_{Y,1}^3 + \{a_1 - a_3 - 3(a_2 - a_3)\} s_{Y,3} + 3(a_2 - a_3) s_{Y,2} s_{Y,1} \\ b &= b_4 s_{Y,1}^3 + \{b_1 - b_2 - 2b_3 + 2b_4\} s_{Y,3} + 3(b_2 + 2b_3 - 3b_4) s_{Y,2} s_{Y,1} \\ c &= c_3 s_{Y,1}^3 + \{c_1 - c_3 - 3(c_2 - c_3)\} s_{Y,3} + 3(c_2 - c_3) s_{Y,2} s_{Y,1} \end{aligned} \tag{3.23}$$

$$a_1 = c_3 = \frac{1}{N(N+2)(N+4)}, a_2 = b_4 = \frac{N+3}{(N-1)N(N+2)(N+4)}, a_3 = \frac{N^2+3N-2}{N(N-1)(N-2)(N+2)(N+4)}, b_1 = c_2 = \frac{3}{N(N+2)(N+4)}, b_2 = \frac{3(N+3)}{(N-1)N(N+2)(N+4)}, b_3 = \frac{N+1}{(N-1)N(N+2)(N+4)}, \text{ and } c_1 = \frac{15}{N(N+2)(N+4)}.$$

### 3.5.3 Proof of Theorem 1

First, we begin by showing that sample correlations correspond to an inner products after removing components in the direction  $e$ . For this extend  $e$  to an orthogonal basis  $e_1 = e, e_2, \dots, e_n$  of  $\mathbb{R}^n$ . Let  $T$  be the  $(n-1) \times n$  matrix that has rows give by  $e_i^t$  for  $i = 2, 3 \dots n$ . Then  $\begin{bmatrix} T^t e \end{bmatrix} \in \mathcal{O}_n$  and hence we have:

$$\begin{aligned} Te &= 0_{n-1} \\ TT^t &= I_{n-1} \\ T^t T &= I_n - ee^t. \end{aligned} \tag{3.24}$$

The following lemma shows that the covariances can be considered as inner products after transformation by  $T$ .

**Lemma 3.** *Let  $X$  and  $Y$  be two data matrices with  $n$  rows. Let  $\text{Cov}(X), \text{Cov}(Y), \text{Cov}(X, Y)$  denote the sample covairance matrix of  $X$ , of  $Y$ , and between  $X$  and  $Y$ , respectively. Then*

$$\begin{aligned} \text{Cov}(X) &= \frac{1}{N}(TX)^t(TX) \\ \text{Cov}(Y) &= \frac{1}{N}(TY)^t(TY) \\ \text{Cov}(X, Y) &= \frac{1}{N}(TX)^t(TY) \end{aligned} \tag{3.25}$$

where  $N = n - 1$ .

*Proof.* We will only show this result for  $\text{Cov}(X, Y)$  since the result for  $\text{Cov}(X)$  and  $\text{Cov}(Y)$  can be obtained by taking  $Y = X$ . Note that  $ee^t A$  gives the matrix of column means that can be used for

centering  $A$ , for any  $A \in \{X, Y\}$ . Hence

$$\begin{aligned}
\text{Cov}(X, Y) &\doteq \frac{1}{N} (X - ee^t X)^t (Y - ee^t Y) \\
&= \frac{1}{N} ((I_n - ee^t)X)^t ((I_n - ee^t)Y) \\
&= \frac{1}{N} (T^t T X)^t (T^t T Y) \\
&= \frac{1}{N} (TX)^t (TT^t) (TY) \\
&= \frac{1}{N} (TX)^t (TY)
\end{aligned}$$

where the last three lines follow using properties of  $T$  in (3.24).  $\square$

In the following, with  $N = n - 1$ , we will work with independent random matrices  $\mathbf{X} \in \mathbb{R}^{N \times p}$  and  $\mathbf{Y} \in \mathbb{R}^{N \times q}$  that have columns on the unit sphere in  $\mathbb{R}^N$  and the columns of the matrix  $\mathbf{X}$  have a jointly spherical symmetric distribution, i.e.

$$U\mathbf{X} \stackrel{d}{=} \mathbf{X} \quad \text{for each } U \in \mathcal{O}_N. \quad (3.26)$$

The following lemma provides the connection of  $(\mathbf{X}, \mathbf{Y})$  with Theorem 1.

**Lemma 4.** *Let  $X \in \mathbb{R}^{n \times p}$  and  $Y \in \mathbb{R}^{n \times q}$  be fixed data matrices and let  $[\tilde{X}, \tilde{Y}] = [\mathbf{Q}'X, Y]$  where  $\mathbf{Q}' \in \mathcal{G}_n$  is chosen according to the uniform measure. With  $N = n - 1$ , define  $\mathbf{X} \in \mathbb{R}^{N \times p}$  to be the matrix obtained by scaling all columns of  $T\tilde{X}$  to have unit norm. Define  $\mathbf{Y} \in \mathbb{R}^{N \times q}$  similarly in terms of  $T\tilde{Y}$ . Then  $\mathbf{X}$  has a spherically symmetric distribution (3.26), and  $\mathbf{X}^t \mathbf{X}$ ,  $\mathbf{Y}^t \mathbf{Y}$  and  $\mathbf{X}^t \mathbf{Y}$  are the sample correlation matrices of  $\tilde{X}$ , of  $\tilde{Y}$ , and between  $\tilde{X}$  and  $\tilde{Y}$ , respectively.*

*Proof.* Note that  $\mathbf{X} = T\tilde{X} \text{diag}[(T\tilde{X})^t(T\tilde{X})]^{-1/2}$  and  $\mathbf{Y} = T\tilde{Y} \text{diag}[(T\tilde{Y})^t(T\tilde{Y})]^{-1/2}$  where  $\text{diag}[A]$  denotes the diagonal matrix obtained by setting all the off-diagonal entries of a square matrix  $A$  to 0. It follows from Lemma 3 with  $X = \tilde{X}$  and  $Y = \tilde{Y}$  that the sample cross-correlation matrix between  $\tilde{X}$  and  $\tilde{Y}$  is given by

$$\begin{aligned}
\text{Cor}(\tilde{X}, \tilde{Y}) &\doteq \text{diag}[\text{Cov}(\tilde{X})]^{-1/2} \text{Cov}(\tilde{X}, \tilde{Y}) \text{diag}[\text{Cov}(\tilde{Y})]^{-1/2} \\
&= \text{diag}[(TX)^t(TX)]^{-1/2} (TX)^t(TY) \text{diag}[(TY)^t(TY)]^{-1/2} \\
&= \mathbf{X}^t \mathbf{Y}.
\end{aligned} \quad (3.27)$$

Similarly, we can show that  $\text{Cor}(\tilde{X}) = \mathbf{X}^t \mathbf{X}$  and  $\text{Cor}(\tilde{Y}) = \mathbf{Y}^t \mathbf{Y}$ . Hence the proof will be complete once we show that  $\mathbf{X}$  satisfies (3.26).

To show  $\mathbf{X}$  satisfies (3.26), let us examine the isomorphism between  $\Psi : \mathcal{G}_n \mapsto \mathcal{O}_{n-1}$  that was alluded to earlier. One may take

$$\Psi(Q) \doteq TQT^t. \quad (3.28)$$

This is indeed a group homomorphism as can be checked using (3.24) that  $\Psi(I_n) = TT^t = I_n$  and

$$\begin{aligned} \Psi(Q)\Psi(P) &= TQ(T^tT)PT^t = TQ(I_n - ee^t)PT^t \\ &= T(Q - ee^t)PT^t = T(QP)T^t = \Psi(PQ) \end{aligned} \quad (3.29)$$

where in the second to last line we have used that  $Qe = e$  and in the last line that  $Te = 0$ . As claimed, it is in fact an isomorphism since  $\Psi(T^tUT + ee^t) = U$  for any  $U \in \mathcal{O}_{n-1}$  and if  $\Psi(Q) = I_{n-1}$  then

$$I_n - ee^t = T^t I_{n-1} T = T^t T Q T^t T = (I_n - ee^t) Q (I_n - ee^t) = Q - ee^t \implies Q = I_n.$$

This means that if  $\mathbf{Q}' \in \mathcal{G}_n$  is distributed according to the Haar measure on  $\mathcal{G}_n$ ,

$$\mathbf{Q}' \stackrel{d}{=} \Psi^{-1}(\mathbf{Q}) = T^t \mathbf{Q} T + ee^t$$

where  $\mathbf{Q}$  is distributed according to the Haar measure on  $\mathcal{O}_{n-1}$ . Hence for any  $U \in \mathcal{O}_{n-1}$

$$UT\mathbf{Q}' \stackrel{d}{=} UT(T^t \mathbf{Q} T + ee^t) = U\mathbf{Q} T + ee^t \stackrel{d}{=} \mathbf{Q} T + ee^t = T(T^t \mathbf{Q} T + ee^t) \stackrel{d}{=} T\mathbf{Q}' \quad (3.30)$$

where the middle equality uses  $U\mathbf{Q} \stackrel{d}{=} \mathbf{Q}$ , which follows since  $\mathbf{Q} \in \mathcal{O}_{n-1}$  is distributed according to the Haar measure. Now with (3.30) we can now show that  $\mathbf{X}$  satisfies (3.26) since for any  $U \in \mathcal{O}_{n-1}$

$$\begin{aligned} U\mathbf{X} &= UT\tilde{X} \text{diag}[(T\tilde{X})^t(T\tilde{X})]^{-1/2} \\ &= UT\mathbf{Q}' X \text{diag}[(T\mathbf{Q}' X)^t(T\mathbf{Q}' X)]^{-1/2} \\ &= UT\mathbf{Q}' X \text{diag}[(UT\mathbf{Q}' X)^t(UT\mathbf{Q}' X)]^{-1/2} \\ &\stackrel{d}{=} T\mathbf{Q}' X \text{diag}[(T\mathbf{Q}' X)^t(T\mathbf{Q}' X)]^{-1/2} = \mathbf{X} \end{aligned} \quad (3.31)$$

where we have used (3.30) to obtain the first inequality in the last line.  $\square$

Apart from the above the above Lemma, the spherical symmetry (3.26) can also be satisfied in other instances, for example if the matrix  $\mathbf{X}$  is obtained by normalizing the columns of a matrix  $\tilde{\mathbf{X}}$  that has i.i.d rows distributed according to a multivariate normal distribution. Hence we will prove the following theorem, which generalizes Theorem 1 with  $N = n - 1$ .

**Theorem 3.** *Let  $\mathbf{X} \in \mathbb{R}^{N \times p}$  and  $\mathbf{Y} \in \mathbb{R}^{N \times p}$  be independent random matrices. Assume that  $\mathbf{X}$  satisfies the spherical symmetry condition, i.e  $U\mathbf{X} \stackrel{d}{=} \mathbf{X}$  for every  $U \in \mathcal{O}_N$ , and that the columns of  $\mathbf{X}$  and  $\mathbf{Y}$  have unit norm in  $\mathbb{R}^N$ . Then*

$$\|\mathbf{R}\|_F^2 = \lambda_{\mathbf{X}}^t (\mathbf{Q} \odot \mathbf{Q}) \lambda_{\mathbf{Y}}^t \quad (3.32)$$

where  $\mathbf{R} = \mathbf{X}^t \mathbf{Y}$ ,  $\lambda_A \in \mathbb{R}_+^N$  is made up of the non-zero eigenvalues of the matrix  $A^t A$  (padded with zeros if necessary) for  $A \in \{\mathbf{X}, \mathbf{Y}\}$ , and  $\mathbf{Q}$  is an random orthogonal matrix distributed according to the Haar measure on  $\mathcal{O}_N$ , independent of  $\lambda_{\mathbf{X}}$  and  $\lambda_{\mathbf{Y}}$ .

The proof of Theorem 3 will proceed by SVD of  $\mathbf{X}$  and  $\mathbf{Y}$  matrices and the following key fact.

**Lemma 5.** *Suppose the matrix  $\mathbf{X} \in \mathbb{R}^{N \times p}$  has spherically symmetric distribution, i.e.  $U\mathbf{X} \stackrel{d}{=} \mathbf{X}$  for each  $U \in \mathcal{O}_N$ , then it has an S.V.D  $\mathbf{X} = \mathbf{Q} \Lambda_{\mathbf{X}} \mathbf{P}^t$  so that  $\mathbf{Q} \in \mathcal{O}_N$  is distributed according to the Haar measure independent of  $\Lambda_{\mathbf{X}}$  and  $\mathbf{P}$ .*

*Proof.* Consider a deterministic SVD procedure that first calculates  $\mathbf{X}^t \mathbf{X} = \mathbf{P} \Lambda_{\mathbf{X}}^2 \mathbf{P}$  and then uses those to calculate  $\mathbf{Q}$ . Let us denote  $\mathbf{Q} = F(\mathbf{X})$ . Such a procedure must then satisfy,  $F(U\mathbf{X}) = U\mathbf{Q}$  since for any  $U \in \mathcal{O}_N$  since  $\mathbf{X}^t \mathbf{X} = (U\mathbf{X}^t)(U\mathbf{X})$ .

We will show that  $\mathbf{Q}$  calculate by such a procedure for  $\mathbf{X}$  is independent of  $\mathbf{X}^t \mathbf{X}$  and distributed according to the Haar measure. To see this, fix any measurable subset  $A \subset \mathbb{R}^{p \times p}$  of  $p \times p$  matrices such that  $\mathbf{P}(\mathbf{X}^t \mathbf{X} \in A) > 0$ , and define the measure  $\mu$  on  $\mathcal{O}_N$  given by

$$\mu_A(B) = \mathbf{P}(\mathbf{Q} \in B \mid \mathbf{X}^t \mathbf{X} \in A) \quad \text{for any measurable subset } B \subseteq \mathcal{O}_N. \quad (3.33)$$

Next, note for any  $U \in \mathcal{O}_N$  that

$$\begin{aligned}
\mu_A(U^{-1}B) &= \frac{P(\mathbf{Q} \in U^{-1}B, \mathbf{X}^t\mathbf{X} \in A)}{P(\mathbf{X}^t\mathbf{X} \in A)} = \frac{P(F(\mathbf{X}) \in U^{-1}B, \mathbf{X}^t\mathbf{X} \in A)}{P(\mathbf{X}^t\mathbf{X} \in A)} \\
&= \frac{P(F(U\mathbf{X}) \in B, (U\mathbf{X})^t(U\mathbf{X}) \in A)}{P(\mathbf{X}^t\mathbf{X} \in A)} \\
&= \frac{P(F(\mathbf{X}) \in B, \mathbf{X}^t\mathbf{X} \in A)}{P(\mathbf{X}^t\mathbf{X} \in A)} = \mu_A(B)
\end{aligned} \tag{3.34}$$

where for the first inequality in the last line, we have used that  $U\mathbf{X} \stackrel{d}{=} \mathbf{X}$ . Since  $\mu_A$  is invariant under left action by any  $U \in \mathcal{O}_N$ ,  $\mu_A$  must be the same Haar measure, regardless of the choice of  $A$ . This shows that  $P(\mathbf{Q} \in B, \mathbf{X}^t\mathbf{X} \in A) = P(\mathbf{Q} \in B)P(\mathbf{X}^t\mathbf{X} \in A)$  for any subsets  $A$  and  $B$  and hence  $\mathbf{Q}$  is distributed according to the Haar measure independent of  $\mathbf{X}^t\mathbf{X}$  (and in turn of  $\Lambda_{\mathbf{X}}$  and  $\mathbf{P}$  which are obtained from  $\mathbf{X}^t\mathbf{X}$ ).  $\square$

*Proof of Theorem 3.* Note that  $\|\mathbf{R}\|_F^2 = \text{tr}(\mathbf{R}^t\mathbf{R}) = \text{tr}(\mathbf{Y}^t\mathbf{X}\mathbf{X}^t\mathbf{Y}) = \text{tr}(\mathbf{Y}\mathbf{Y}^t\mathbf{X}\mathbf{X}^t)$ . Let  $\mathbf{Y} = U\Lambda_{\mathbf{Y}}V^t$  be any SVD and let  $\mathbf{X} = \tilde{\mathbf{Q}}\Lambda_{\mathbf{X}}\mathbf{P}^t$  be the SVD from Lemma 5 so that  $\tilde{\mathbf{Q}} \in \mathcal{O}_N$  is distributed according to the Haar measure. Then

$$\begin{aligned}
\|\mathbf{R}\|_F^2 &= \text{tr}(\mathbf{Y}\mathbf{Y}^t\mathbf{X}\mathbf{X}^t) = \text{tr}(U\Lambda_{\mathbf{Y}}^2U^t\tilde{\mathbf{Q}}\Lambda_{\mathbf{X}}\tilde{\mathbf{Q}}^t) \\
&= \text{tr}(\Lambda_{\mathbf{Y}}^2U^t\tilde{\mathbf{Q}}\Lambda_{\mathbf{X}}^2\tilde{\mathbf{Q}}^tU) = \text{tr}(\Lambda_{\mathbf{Y}}^2\mathbf{Q}\Lambda_{\mathbf{X}}^2\mathbf{Q}^t) = \langle \Lambda_{\mathbf{X}}^2\mathbf{Q}, \mathbf{Q}\Lambda_{\mathbf{Y}}^2 \rangle
\end{aligned} \tag{3.35}$$

where  $\mathbf{Q} = U^t\tilde{\mathbf{Q}}$  and  $\langle A, B \rangle = \text{tr}(AB^t)$  is the usual inner product where the matrices are considered as vectors. Since  $\mathbf{X}$  is independent of  $\mathbf{Y}$ ,  $\mathbf{Q}$  is distributed uniformly on  $\mathcal{O}_N$ , independent of both  $\Lambda_{\mathbf{X}}^2$  and  $\Lambda_{\mathbf{Y}}^2$ . Finally note that  $\Lambda_{\mathbf{X}}^2$  and  $\Lambda_{\mathbf{Y}}^2$  are diagonal  $N \times N$  matrices with diagonal entries given by the vectors  $\lambda_{\mathbf{X}}$  and  $\lambda_{\mathbf{Y}}$  – the non-zero eigenvectors of  $\mathbf{X}^t\mathbf{X}$  and  $\mathbf{Y}^t\mathbf{Y}$  padded with zeros as necessary. Hence if  $\mathbf{Q} = (Q_{ij})$  then

$$\|\mathbf{R}\|_F^2 = \sum_{i=1}^N \sum_{j=1}^N \lambda_{\mathbf{X},i} Q_{i,j}^2 \lambda_{\mathbf{Y},j} = \lambda_{\mathbf{X}}^t (\mathbf{Q} \odot \mathbf{Q}) \lambda_{\mathbf{Y}} \tag{3.36}$$

$\square$

### 3.5.4 Proof of Theorem 2

Let  $\mathbf{Q} \in \mathcal{O}_N$  be distributed according to the Haar measure, and denote the rows of  $\mathbf{Q}$  by  $\mathbf{q}_1, \dots, \mathbf{q}_N \in \mathbb{R}^N$  and further entries of the matrix as  $\mathbf{Q} = (q_{i,j})$ . We use  $\langle \cdot, \cdot \rangle$  to denote the standard dot product in  $\mathbb{R}^N$ , and let  $e_1 \dots e_N$  be the usual orthonormal basis (i.e.  $e_i = (0, 0, \dots, 1, \dots, 0)$ ).

In the following, we use  $\lambda, \mu \in \mathbb{R}^N$  instead of  $\lambda_X, \lambda_Y$  as defined in Theorem 1. We would like to calculate upto three moments of  $T$  using (3.21). For this we need to be able to compute moments of  $\mathbf{Q}$  of the form  $\mathbf{E}[q_{i_1, j_1}^2 q_{i_2, j_2}^2 q_{i_3, j_3}^3] = \mathbf{E}[\langle \mathbf{q}_{i_1}, e_{j_1} \rangle^2 \langle \mathbf{q}_{i_2}, e_{j_2} \rangle^2 \langle \mathbf{q}_{i_3}, e_{j_3} \rangle^2]$ . We will now demonstrate how to do this.

First we start with some elementary properties of the Haar measure on  $\mathcal{O}_N$  (87).

**Proposition 4.**  $\mathbf{Q} \in \mathcal{O}_N$  distributed according to the Haar measure, satisfies the following properties

1. For any  $W \in \mathcal{O}_N$ ,  $W\mathbf{Q} \stackrel{d}{=} \mathbf{Q}W \stackrel{d}{=} W$
2.  $\mathbf{Q} \stackrel{d}{=} \mathbf{Q}^t$
3. For any  $i \geq 1$ , conditioned on  $\mathbf{q}_1 \dots \mathbf{q}_{i-1}$ ,  $\mathbf{q}_i$  is distributed uniformly on the unit sphere in  $\text{span}\{\mathbf{q}_1 \dots \mathbf{q}_{i-1}\}^\perp$

By taking  $W$  to be a permutation matrix in (1), note that the distribution of the rows (or columns) of  $\mathbf{Q}$  is invariant under permutations. As seen in (3), the Haar measure on  $\mathcal{O}_N$  is closely related to the uniform distribution on the unit sphere. Hence the following result will be a key tool.

**Lemma 6.** *Let  $\mathbf{q}$  be distributed uniformly on the unit sphere of a  $k$  dimensional subspace  $V$  of some Hilbert space. Let  $v_1, \dots, v_k$  be an orthonormal basis for  $V$ . Then  $(\langle \mathbf{q}, v_i \rangle^2)_{\{i \leq k\}}$  is distributed as a  $k$ -dimensional Dirichlet distribution with parameter  $\alpha = (1/2, 1/2, \dots, 1/2) \in \mathbb{R}^k$ .*

*Proof.* Use that  $\mathbf{q} \stackrel{d}{=} (\sum_{i=1}^k Z_i v_i) / \sqrt{\sum_{i=1}^k Z_i^2}$ , where  $Z_i \sim N(0, 1)$  are i.i.d, and then use the representation of Dirichlet in terms of independent Gamma's.  $\square$

### 3.5.5 Calculating moments of $Q \odot Q$

Now we will show how to calculate the moments of  $\mathbf{Q} \odot \mathbf{Q}$  using the results mentioned above. First note by Lemma 6 that  $(\langle \mathbf{q}_1, e_i \rangle^2)_{i=1}^N \stackrel{d}{=} (\pi_{i:N})_{i=1}^N \sim \text{Dir}(1/2, 1/2 \dots 1/2) \in \mathbb{R}^N$ . Hence

$$\begin{aligned} \mathbf{E} \langle \mathbf{q}_1, e_1 \rangle^2 &= \mathbf{E} \pi_{1:N} \\ &= \frac{1}{N} \end{aligned}$$

Arbitrary moments of Dirichlet are also easy to calculate. Hence this approach could also be used to calculate  $\mathbf{E} \langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_1, e_2 \rangle^4$ .

Now let us try to calculate a more complicated term like  $\mathbf{E} \langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_2, e_2 \rangle^2$ . We will use that, when conditioned on  $\mathbf{q}_1$ ,  $\mathbf{q}_2$  is distributed uniformly on the space perpendicular to  $\mathbf{q}_1$ . Hence  $\langle \mathbf{q}_2, e_2 \rangle = \langle \mathbf{q}_2, e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1 \rangle = \langle \mathbf{q}_2, \tilde{e}_2 \rangle \|e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1\|$ , where  $\tilde{e}_2 = \frac{e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1}{\|e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1\|}$ . Using this, we get:

$$\begin{aligned} \mathbf{E} \langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_2, e_2 \rangle^2 &= \mathbf{E} \left( \langle \mathbf{q}_1, e_1 \rangle^2 \|e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1\|^2 \mathbf{E} \{ \langle \mathbf{q}_2, \tilde{e}_2 \rangle^2 \mid \mathbf{q}_1 \} \right) \\ &= \mathbf{E} \left( \langle \mathbf{q}_1, e_1 \rangle^2 \|e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1\|^2 \mathbf{E} \{ \pi_{1:N-1} \} \right) \\ &= \mathbf{E} \left( \langle \mathbf{q}_1, e_1 \rangle^2 \|e_2 - \langle \mathbf{q}_1, e_2 \rangle \mathbf{q}_1\|^2 \frac{1}{N-1} \right) \\ &= \frac{1}{N-1} \mathbf{E} (\langle \mathbf{q}_1, e_1 \rangle^2 (1 - \langle \mathbf{q}_1, e_2 \rangle^2)) \\ &= \frac{1}{N-1} \mathbf{E} (\pi_{1:N} - \pi_{2:N} \pi_{1:N}) \\ &= \frac{1}{N-1} \left( \frac{1}{N} - \frac{1}{N(N+2)} \right) \\ &= \frac{N+1}{(N-1)N(N+2)} \end{aligned} \tag{3.37}$$

This technique can be similarly used to calculate  $\mathbf{E} \langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_2, e_2 \rangle^2 \langle \mathbf{q}_3, e_3 \rangle^2$ , by first conditioning on  $\mathbf{q}_1, \mathbf{q}_2$  and eliminating the  $\mathbf{q}_3$  term, and then eliminating the  $\mathbf{q}_2$  term by conditioning on  $\mathbf{q}_1$ .

### 3.5.6 Calculating moments of $T$

Once we can calculate the moments of  $\mathbf{Q} \odot \mathbf{Q}$ , in principle it should be possible to calculate arbitrary moments of  $T$ . However the expressions would be in terms of many summations. Since



the rows and columns of  $\mathbf{Q}$  are symmetric under permutations, the formula for moments of  $T$  can be simplified. For instance, since  $\forall i, j \ \mathbf{E}q_{ij}^2 = \mathbf{E}q_{11}^2$ , we have:

$$\mathbf{E}T = \left(\sum_{i=1}^N \lambda_i\right) \left(\sum_{j=1}^N \mu_j\right) \mathbf{E}q_{11}^2 = \left(\sum_{i=1}^N \lambda_i\right) \left(\sum_{j=1}^N \mu_j\right) / N \quad (3.38)$$

### 3.5.6.1 The second moment of $T$

To use this pattern for higher moments, it will be useful to proceed systematically. First write  $T$  in a concise notation:

$$T = \sum_{i=1}^N \lambda_i \langle \mathbf{q}_i^{\odot 2}, \mu \rangle$$

Where  $\mathbf{q}^{\odot 2}$  refers to the vector obtained by squaring each coordinate of  $\mathbf{q}_i$ . Hence

$$\mathbf{E}T^2 = \sum_{i_1, i_2=1}^N \lambda_{i_1} \lambda_{i_2} \mathbf{E} \langle \mathbf{q}_{i_1}^{\odot 2}, \mu \rangle \langle \mathbf{q}_{i_2}^{\odot 2}, \mu \rangle \quad (3.39)$$

Let  $a = \mathbf{E} \langle \mathbf{q}_1^{\odot 2}, \mu \rangle \langle \mathbf{q}_2^{\odot 2}, \mu \rangle$  and  $b = \mathbf{E} \langle \mathbf{q}_1^{\odot 2}, \mu \rangle^2$ . Then (3.39) becomes

$$\begin{aligned} \mathbf{E}T^2 &= \sum_{i_1, i_2=1}^N \lambda_{i_1} \lambda_{i_2} (a \mathbb{I}_{\{i_1 \neq i_2\}} + b \mathbb{I}_{\{i_1 = i_2\}}) \\ &= \sum_{i_1, i_2=1}^N \lambda_{i_1} \lambda_{i_2} (a(1 - \mathbb{I}_{\{i_1 = i_2\}}) + b \mathbb{I}_{\{i_1 = i_2\}}) \\ &= \sum_{i_1, i_2=1}^N \lambda_{i_1} \lambda_{i_2} [a + (b - a) \mathbb{I}_{\{i_1 = i_2\}}] \\ &= a \left( \sum_{i=1}^N \lambda_i \right)^2 + (b - a) \left( \sum_{i=1}^N \lambda_i^2 \right) \end{aligned} \quad (3.40)$$

Now to find  $a$ , let  $a_1 = \mathbf{E}\langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_2, e_1 \rangle^2$  and  $a_2 = \mathbf{E}\langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_2, e_2 \rangle^2$ . Then:

$$\begin{aligned}
a &= \sum_{j_1, j_2=1}^N \mu_{j_1} \mu_{j_2} \mathbf{E}(q_{1,j_1}^2 q_{2,j_2}^2) \\
&= \sum_{j_1, j_2=1}^N \mu_{j_1} \mu_{j_2} (a_1 \mathbb{I}_{\{j_1=j_2\}} + a_2 \mathbb{I}_{\{j_1 \neq j_2\}}) \\
&= (a_1 - a_2) \left( \sum_{j=1}^{n-1} \mu_j^2 \right) + a_2 \left( \sum_{j=1}^{n-1} \mu_j \right)^2
\end{aligned}$$

We have already shown in (3.37) that  $a_2 = \frac{N+1}{(n-1)N(N+2)}$ . We can find the value of  $a_1$  similarly, but instead we will use that  $\mathbf{Q} \stackrel{d}{=} \mathbf{Q}^t$ . This would show  $a_1 = \mathbf{E}\langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_1, e_2 \rangle^2 = \mathbf{E}\pi_{1:N}\pi_{2:N} = \frac{1}{N(N+2)}$ .

Now to find  $b$ , let  $b_1 = \mathbf{E}\langle \mathbf{q}_1, e_1 \rangle^4$  and  $b_2 = \mathbf{E}(\langle \mathbf{q}_1, e_1 \rangle^2 \langle \mathbf{q}_1, e_2 \rangle^2)$ . Then

$$\begin{aligned}
b &= \sum_{j_1, j_2=1}^N \mu_{j_1} \mu_{j_2} \mathbf{E}(q_{1,j_1}^2 q_{1,j_2}^2) \\
&= \sum_{j_1, j_2=1}^N \mu_{j_1} \mu_{j_2} \mathbf{E}\{(b_1 - b_2) \mathbb{I}_{\{j_1=j_2\}} + b_2\} \\
&= (b_1 - b_2) \left( \sum_{j=1}^N \mu_j^2 \right) + b_2 \left( \sum_{j=1}^{n-1} \mu_j \right)^2
\end{aligned}$$

Here  $b_1 = \mathbf{E}\pi_{1:N}^2 = \frac{3}{N(N+2)}$  and  $b_2 = \mathbf{E}\pi_{1:N}\pi_{2:N} = \frac{1}{N(N+2)}$ . Combining all of this provides a formula for  $\mathbf{E}T^2$ .

### 3.5.6.2 The third moment of $T$

The third moment of  $T$  can also be obtained using the same techniques as above. In particular start with

$$\mathbf{E}T^3 = \sum_{i_1, i_2, i_3} \lambda_{i_1} \lambda_{i_2} \lambda_{i_3} \mathbf{E}\langle \mathbf{q}_{i_1}^{\odot 2}, \mu \rangle \langle \mathbf{q}_{i_2}^{\odot 2}, \mu \rangle \langle \mathbf{q}_{i_3}^{\odot 2}, \mu \rangle \quad (3.41)$$

where  $a \doteq \mathbf{E}\langle \mathbf{q}_1^{\odot 2}, \mu \rangle \langle \mathbf{q}_2^{\odot 2}, \mu \rangle \langle \mathbf{q}_3^{\odot 2}, \mu \rangle$ ,  $b \doteq \mathbf{E}\langle \mathbf{q}_1^{\odot 2}, \mu \rangle^2 \langle \mathbf{q}_2^{\odot 2}, \mu \rangle$ , and  $c \doteq \mathbf{E}\langle \mathbf{q}_1^{\odot 2}, \mu \rangle^3$ , and we proceeded to simplify  $T, a, b, c$  using the same method as in the  $\mathbf{E}T^2$ . For brevity, we omit the

calculation details and provide the final formula for the third moment. In the following  $N = n - 1$  and The summations are take over indices in the set  $\{1 \dots N\}$ .

$$\mathbf{E}T^3 = a \left( \sum_i \lambda_i \right)^3 + \{c - a - 3(b - a)\} \left( \sum_i \lambda_i^3 \right) + 3(b - a) \left( \sum_i \lambda_i^2 \right) \left( \sum_i \lambda_i \right)$$

where

$$\begin{aligned} a &= a_3 \left( \sum_j \mu_j \right)^3 + \{a_1 - a_3 - 3(a_2 - a_3)\} \left( \sum_j \mu_j^3 \right) + 3(a_2 - a_3) \left( \sum_j \mu_j^2 \right) \left( \sum_j \mu_j \right) \\ b &= b_4 \left( \sum_j \mu_j \right)^3 + \{b_1 - b_2 - 2b_3 + 2b_4\} \left( \sum_j \mu_j^3 \right) + 3(b_2 + 2b_3 - 3b_4) \left( \sum_j \mu_j^2 \right) \left( \sum_j \mu_j \right) \\ c &= c_3 \left( \sum_j \mu_j \right)^3 + \{c_1 - c_3 - 3(c_2 - c_3)\} \left( \sum_j \mu_j^3 \right) + 3(c_2 - c_3) \left( \sum_j \mu_j^2 \right) \left( \sum_j \mu_j \right) \end{aligned}$$

where  $a_1 \doteq \mathbf{E}(q_{11}q_{12}q_{13})^2 = \frac{1}{N(N+2)(N+4)}$ ,  $a_2 \doteq \mathbf{E}(q_{11}q_{12}q_{23})^2 = \frac{N+3}{(N-1)N(N+2)(N+4)}$ ,  
 $a_3 \doteq \mathbf{E}(q_{11}q_{22}q_{33})^2 = \frac{N^2+3N-2}{N(N-1)(N-2)(N+2)(N+4)}$ ,  $b_1 \doteq \mathbf{E}(q_{11}q_{11}q_{12})^2 = \frac{3}{N(N+2)(N+4)}$ ,  $b_2 \doteq \mathbf{E}(q_{11}q_{11}q_{22})^2 = \frac{3(N+3)}{(N-1)N(N+2)(N+4)}$ ,  $b_3 \doteq \mathbf{E}(q_{11}q_{21}q_{12})^2 = \frac{N+1}{(N-1)N(N+2)(N+4)}$ ,  $b_4 \doteq \mathbf{E}(q_{11}q_{21}q_{32})^2 = \frac{N+3}{(N-1)N(N+2)(N+4)}$ ,  $c_1 \doteq \mathbf{E}q_{11}^6 = \frac{15}{N(N+2)(N+4)}$ ,  $c_2 \doteq \mathbf{E}(q_{11}^2q_{12})^2 = \frac{3}{N(N+2)(N+4)}$ , and  $c_3 \doteq \mathbf{E}(q_{11}q_{12}q_{13})^2 = \frac{1}{N(N+2)(N+4)}$ .

## CHAPTER 4

### Data analysis

In this chapter, some of the methods for bimodule detection that were discussed previously will be applied to an artificial and real data set. First, in Section 4.1, we simulate a large dataset in which several true bimodules are planted, and we then measure the effectiveness of the various methods at recovering these true bimodules.

Next, in Section 4.2, we apply the bimodule detection methods to an eQTL dataset from the GTEx consortium. In particular, for the SNP-gene bimodules produced by BSP, we examine a variety of descriptive and biological metrics, including network structure, genomic locations, enrichment for known gene sets from the gene-ontology database, and comparisons with, and potential extensions of, standard eQTL analysis.

Although not discussed in this dissertation, a further application of BSP to a climate data-set consisting of inter-annual temperature and precipitation measurements in North America can be found in (41).

#### 4.1 Simulation study

To assess the effectiveness of BSP, we carried out a simulation study in which a variety of true bimodules of different strengths and sizes were present in the underlying distribution of the samples. In this section, we provide an overview of the study, and an assessment of the results from BSP and competing methods CONDOR and sCCA (see Sections 2.2.2 and 2.2.1).

Simulation studies incorporating fewer than ten embedded bimodules have been conducted for methods based on sCCA (125, 102, 130) and graphical models (33, 32). However, existing studies are relatively simple, and do not emphasize the network structure of many applications. In order to emulate the complexity of eQTL analysis and similar applications, we designed a simulation study in which  $K = 500$  bimodules of various strengths, sizes, network structures, and intra-

correlations were planted in a single large dataset. The planted bimodules were then connected by confounding edges to make their recovery more challenging. We emphasize that BSP is not based on an underlying generative model: the model used in the simulation study is for assessment purposes only.

#### 4.1.1 Details of the simulated data

We generated a single large dataset having  $n = 200$  samples and two measurement types, with  $p = 100,000$  and  $q = 20,000$  features, respectively. The number of features is of the same order of magnitude as in the eQTL dataset considered in Section 4.2. Following the notation in Section 3.1, we denote the two types of features by index sets  $S = \{s_1, s_2 \dots s_p\}$  and  $T = \{t_1, t_2 \dots t_q\}$ . For each individual, the joint  $p + q$  dimensional measurement vector is independently drawn from a multivariate normal distribution with mean  $0 \in \mathbb{R}^{p+q}$  and  $(p + q) \times (p + q)$  covariance matrix  $\Sigma$ . The covariance matrix  $\Sigma$  is designed so that it has  $K = 500$  true bimodules of various sizes, network structures, signal strengths and intra-correlations. As it is difficult to generate structured covariance matrices while maintaining non-negative definiteness, we instead specify a generative model for the  $p + q$  dimensional random row vector  $(X, Y) \sim \mathcal{N}_{p+q}(0, \Sigma)$ .

We now describe how we generated the block cross-correlation signal between the two feature types, representing observed bimodules. To begin, we partitioned the first-half of the  $S$ -indices  $\{s_1, \dots, s_{\lceil p/2 \rceil}\}$  into  $K$  disjoint subsets  $A_1, A_2, \dots, A_K$  with sizes chosen according to a Dirichlet distribution with parameter  $(1, 1, \dots, 1) \in \mathbb{R}^K$ . In the same way, we generated a Dirichlet partition  $B_1, B_2, \dots, B_K$  of the first-half of  $T$  indices  $\{t_1, \dots, t_{\lceil q/2 \rceil}\}$  independent of the previous partition. The feature-set pairs  $(A_i, B_i)$  constitute the true bimodules, while the features in second-half of the  $S$ - and  $T$ -indices are not part of true bimodules. Next, the random sub-vectors  $(X_{A_i}, Y_{B_i})$  corresponding to the true bimodules were generated independently for each  $i \in [K]$  using a graph based regression model described below.

Let  $(A, B)$  be a feature set pair, and suppose that  $\rho \in [0, 1)$  and  $\sigma^2 > 0$  are given. Let  $D \in \{0, 1\}^{|A| \times |B|}$  be a binary matrix, which we regard as the adjacency matrix of a connected bipartite network with vertex set  $A \cup B$ . Then the random row-vector  $(X_A^t, Y_B^t)$  is generated as follows:

$$X_A \sim \mathcal{N}_{|A|}(0, (1 - \rho)I + \rho U) \quad \text{and} \quad Y_B = D^t X_A + \epsilon, \quad (4.1)$$

where  $\epsilon \sim \mathcal{N}_{|B|}(0, \sigma^2 I)$  and  $U$  is a matrix of all ones. To understand the bimodule signal produced by this model, note that  $\rho$  governs the intra-correlation between features in  $A$  and that for any  $t \in B$ , the variable  $Y_t$  is influenced by features  $X_s$  such that  $(s, t)$  is an edge in adjacency matrix  $D$ . For each of the true bimodules  $(A_i, B_i)$  in the simulation, we independently chose parameters  $\rho_i$ ,  $\sigma_i^2$ , and  $D_i$  to produce a variety of behaviors while maintaining the inherent constraints between them. We provide more details in Section 4.1.5.

Among features that are not part of bimodules, features  $X_{s_j}$  with  $j > \lceil p/2 \rceil$  are independent  $\mathcal{N}(0, 1)$  noise variables and features  $Y_{t_r}$  with  $r > \lceil q/2 \rceil$  are either noise variables (generated independently as  $\mathcal{N}(0, 1)$ ) or they are bridge variables that connect two true bimodules. In more detail, for every pair of distinct bimodules  $(A_k, B_k)$  and  $(A_l, B_l)$  with  $1 \leq k < l \leq K$ , with probability  $\frac{1.5}{K}$ , we connect the two bimodules by selecting at random (and without replacement) an index  $r > \lceil q/2 \rceil$  and making it a bridge variable by defining

$$Y_{t_r} = X_s + X_{s'} + \epsilon \quad \text{with } \epsilon \sim N(0, \sigma_r^2), \quad (4.2)$$

for a randomly chosen  $s \in A_k$  and  $s' \in A_l$ . The noise variance  $\sigma_r^2$  in (4.2) is chosen so that the correlation strength between  $Y_{t_r}$  and  $X_s$  (and  $X_{s'}$ ) is equal to the average strength of the bimodules that are being connected. If  $Y_{t_r}$  is not a bridge variable, it is taken to be noise.

Prior to the addition of bridge variables, the connected components of the population cross-correlation network are just the bimodules  $(A_k, B_k)$ . Once bridge variables have been added, the population cross-correlation network will have a so-called giant connected component comprising a substantial portion of the underlying index space  $S \times T$ . While theoretical support for the presence of giant component in our simulation model comes from the study of Erdős-Renyi random graphs (19), such components have also been observed in empirical eQTL networks (47, 107). Although the giant component is itself a stable population bimodule, since we only add a small number (348) of bridge variables, the majority of the cross-correlation signal is in the more densely connected sets  $(A_k, B_k)$ , which we continue to refer to as the *true* bimodules.

### 4.1.2 Running BSP and related methods

We applied BSP to the simulated data using the false discovery parameter  $\alpha = 0.01$ , which was selected to keep the edge-error estimates under 0.05 (see Section 3.4.3). The search was initialized from singletons consisting of all the features in  $T$  and 1% of the features in  $S$ , chosen at random. In what follows, feature-set pairs identified by BSP (or some other method, when clear from context) will be referred to as *detected* bimodules. BSP detected 319 unique bimodules while the effective number (see Section 3.4.4) of detected bimodules was 301.5.

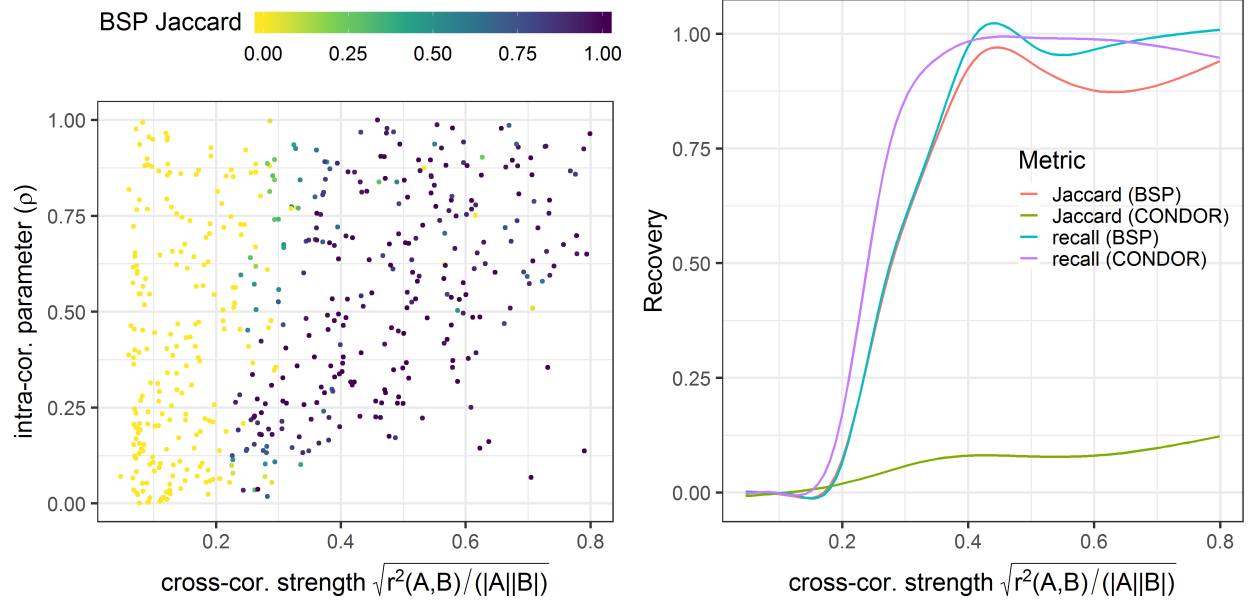
To obtain bimodules via CONDOR (107), we applied Matrix-eQTL (114) to the simulated dataset with  $S$  considered as the set of SNPs and  $T$  considered as the set of genes, to extract feature pairs  $(s, t) \in S \times T$  with q-value less than  $\alpha = 0.05$ . Next, we formed a bipartite graph on the vertex set  $S \cup T$  with edges given by the significant feature pairs found in the previous step. The largest connected component of this graph, made up of 28,876 features from  $S$  and 6,455 features from  $T$ , was passed through a bipartite community detection software (106) which partitioned the nodes of the sub-graph into 112 bimodules.

We applied the sCCA software (129) (based on (130)) to the simulated data to find 100 bimodules. More precisely, for various penalty parameters  $\lambda \in [0, 1]$ , we ran sCCA to find 100 canonical covariate pairs with the  $\ell_1$  norm constraint of  $\lambda\sqrt{p}$  and  $\lambda\sqrt{q}$  on the coefficients of the linear combinations corresponding to  $S$  and  $T$  respectively. Initially, we considered  $\lambda = 0.233$ , chosen by the permutation based procedure provided with the software. However the resulting bimodules were large and had high edge-error (see Section 4.3.1). Based on a rough grid search, we then ran the procedure with each value  $\lambda \in \{.01, .02, .03, .04, .06\}$  to obtain smaller bimodules.

### 4.1.3 Comparing performance of the methods

In the simulation study described above, we measure the recovery of a true bimodule  $(A_t, B_t)$  by a detected bimodule  $(A_d, B_d)$  using the two metrics:

$$\text{recall} = \frac{|A_t \cap A_d| |B_t \cap B_d|}{|A_t| |B_t|} \quad \text{and} \quad \text{Jaccard} = \frac{|A_t \cap A_d| |B_t \cap B_d|}{|A_t \times B_t \cup A_d \times B_d|}.$$



**Figure 4.1:** Recovery of true bimodules. Left: dependence of cross-correlation strength and intra-correlation parameter of true bimodules on BSP Jaccard. Right: the averaged recovery curves (recall and Jaccard) for true bimodules under CONDOR and BSP.

Recall captures how well the true bimodule is *contained* inside the detected bimodule, while Jaccard measures how well the two bimodules *match*. When assessing the recovery of a true bimodule under a collection of detected bimodules (like the output of BSP), we choose the detected bimodule with the best recall or Jaccard, depending on the metric under consideration.

As shown in Figure 4.1, the BSP Jaccard for true bimodules was influenced primarily by the cross-correlation strength  $\sqrt{\frac{r^2(A,B)}{|A||B|}}$  of the true bimodule, though the intra-correlation parameter  $\rho$  used in the simulation (4.1) was also seen to have an effect (Figure 4.1, left). Most bimodules with cross-correlation strength above 0.4 were completely recovered, while those with strength below 0.2 were not recovered. For strengths between 0.2 to 0.4, there was a variation in Jaccard, with smaller Jaccard for bimodules having larger values of  $\rho$  (Figure 4.1, left). The effect of  $\rho$  on Jaccard was expected since BSP accounts for the intra-correlation among features of the same type.

The intra-correlation parameter  $\rho$  did not have significant effect on CONDOR Jaccard, since the method does not account for intra-correlations. Hence, here we only consider the effects of the cross-correlation strength of true bimodules on CONDOR Jaccard (Figure 4.1, green curve on the right). Regardless of the cross-correlation strength, CONDOR Jaccard remained low. This was because CONDOR bimodules often overlapped multiple true bimodules; indeed, 102 of the



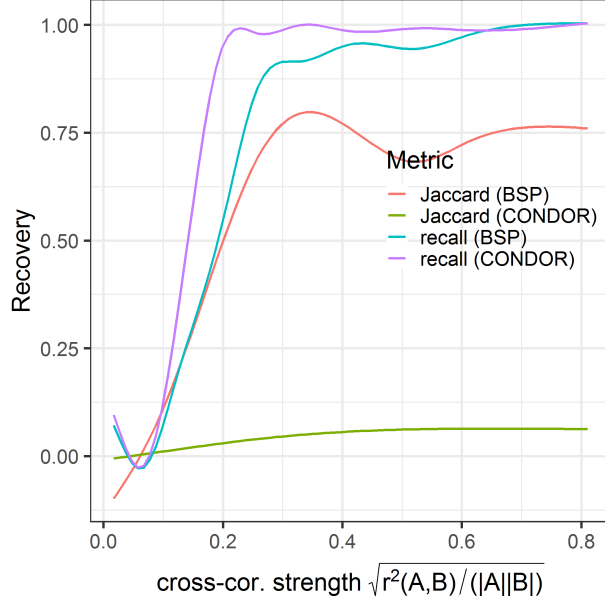
112 CONDOR bimodules overlapped with two or more (up to 19) true bimodules, compared with only 21 of the 319 BSP bimodules. However, the results for CONDOR recall (Figure 4.1, purple curve on the right) show that most true bimodules with significant cross-correlation strengths were contained inside some CONDOR bimodule.

To assess the false discoveries in detected bimodules, we measured the *edge-error* of detected bimodules. The edge-error is the fraction of the essential-edges (see (3.13) in Section 3.3.4) of a detected bimodule that are not part of the simulation model, that is, edges not contained in any true bimodule and not in the set of bridge edges. The average edge-error for BSP bimodules was 0.03, and 90% of the detected bimodules had edge-error under 0.05. In contrast, the average edge-error for CONDOR bimodules was 0.08, and 90% of the detected bimodules had edge-error under 0.14. The larger edge-error among CONDOR bimodules may have arisen because the method does not account for intra-correlations.

Concerning sCCA, the sizes of the detected bimodules were at least an order of magnitude larger than sizes of the true bimodules when  $\lambda$  exceeded 0.04 (Figure 4.6). Thus we only considered  $\lambda \leq 0.04$ . For  $\lambda = 0.03$  and 0.04, the detected bimodules had large edge-error (average error 0.47 and 0.65, respectively), while for  $\lambda = 0.01$  and 0.02 the true bimodules had poor recall (95% of the true bimodules had recall below 0.02 and 0.23, respectively). Further details of the results are given in Section 4.3.1. A potential shortcoming of our application of sCCA was that we chose the same penalty parameter  $\lambda$  for each of the 100 bimodules. We expect that the results of sCCA would improve if one chose a different penalty parameter for each bimodule. However Witten et al. (130) do not provide explicit guidelines to choose different penalty parameters for each component (bimodule), and directly doing a permutation-based grid search each time would be exceedingly slow.

#### 4.1.4 Performance of BSP and CONDOR on increasing sample size

We also studied the performance of BSP and CONDOR on a simulation study with a larger sample size of  $n = 600$ . As expected, both methods were able to recall bimodules with lower cross-correlation strengths than earlier. However, both BSP and CONDOR had lower Jaccard than in the  $n = 200$  simulation. We will discuss details about this analysis now, but see Chapter 8 (Section 8.1.1) for a discussion on why such behavior may be expected.



**Figure 4.2:** Average recall and Jaccard for true bimodules in the simulation with 600 samples.

We generated a new artificial data-set using the procedure in Section 4.1 with  $n = 600$ , and reran BSP and CONDOR with the same parameters as earlier. The average edge-error for BSP and CONDOR was 0.05 and 0.10 respectively. As seen in Figure 4.2, BSP and CONDOR both recall most bimodules with cross-correlation strength above 0.3, however Jaccard for BSP and CONDOR has degraded. This can be explained by noting that 25% of BSP bimodules now overlapped with two or more true bimodules compared to 6% when  $n = 200$ .

We now conclude the discussion on the simulation study with some details about the generative model used earlier.

#### 4.1.5 Choice for parameters used in the simulation

As described in Section 4.1.1, given  $\rho, \sigma \in [0, 1]$  and a binary adjacency matrix  $D \in \{0, 1\}^{|A| \times |B|}$  representing a connected bipartite graph on vertices  $A$  and  $B$  (called the regressor-graph), the random row vector  $(X_A^t, Y_A^t)$  of variables for a bimodule  $(A, B)$  can be simulated as

$$X_A \sim \mathcal{N}_{|A|}(0, (1 - \rho)I + \rho U) \quad \text{and} \quad Y_B = D^t X_A + \epsilon, \quad (4.3)$$

where  $U$  is the matrix of all ones and  $\epsilon \sim \mathcal{N}_{|B|}(0, \sigma^2 I)$ . The parameters  $\rho, \sigma \in [0, 1]$  and  $D$  appearing in (4.3) are chosen independently for each true bimodule  $(A, B)$  as follows:

1. Choose a constant  $\beta \in [0, 1]$  uniformly at random. With  $d \doteq \lceil \beta |A| \rceil$ , let  $D$  be the adjacency matrix of the  $d$ -regular bipartite connected graph on vertex sets  $A$  and  $B$  formed by independently connecting each vertex  $t \in B$  to  $d$  randomly chosen vertices from  $A$ . If the resulting graph is not connected, set  $\beta$  to  $\beta + \Delta\beta$  where  $\Delta\beta = 0.1$  and repeat the previous step till the resulting bipartite graph is connected.
2. Randomly choose  $\rho \in [0, 1]$  and  $\eta \in [0, .8]$  subject to the constraint  $\delta \doteq 1 + \rho(d - 1) \geq \eta^2 d$ . We satisfy this constraint by first uniformly generating  $\rho$  and then generating  $\eta$  uniformly from  $[0, \min(\sqrt{\delta d^{-1}}, .8)]$ .
3. Finally let  $\sigma = \frac{\sqrt{\delta(\delta - \eta^2 d)}}{\eta}$ .

The constants  $(\rho, \beta, \eta)$  in the above procedure have the following intuitive role:  $\rho$  is the intra-correlation between any two features from the set  $A$ ,  $\beta \in [0, 1]$  controls the edge density of the regressor-graph  $D$ , and  $\eta$  is the cross-correlation between features from  $B$  and adjacent features from  $A$  in the regressor-graph. The following Lemma shows that our choice of parameters indeed results in population cross-correlation of  $\eta$  between features connected by the regressor-graph:

**Lemma 7.** *Fix  $\rho, \eta \in [0, 1]$ ,  $a, b \in \mathbb{N}$  and  $d \in \{1, 2, \dots, a\}$  so that  $\delta \doteq 1 + \rho(d - 1) \geq \eta^2 d$ . Suppose  $\mathbf{X} \in \mathbb{R}^a$  is a random vector with covariance matrix  $\text{Cov}(\mathbf{X}) = \rho U_a + (1 - \rho)I_a$ , where  $U_a \in \mathbb{R}^{a \times a}$  is the matrix of all ones and  $I_a \in \mathbb{R}^{a \times a}$  is the identity matrix. Next suppose  $D$  is a  $\{0, 1\}$  valued  $a \times b$  dimensional matrix that has exactly  $d$  ones in each column. Finally let  $\sigma = \sqrt{\delta(\delta - \eta^2 d)}/\eta$  and suppose the  $b$ -dimensional random vector  $\mathbf{Y}$  is given by*

$$\mathbf{Y} = D^t \mathbf{X} + \epsilon$$

where  $\epsilon$  is another  $b$ -dimensional random vector independent of  $\mathbf{X}$  with  $\text{Cov}(\epsilon) = \sigma^2 I_b$ . Then

$$\text{Cor}(\mathbf{X}, \mathbf{Y}) \odot D = \eta D \tag{4.4}$$

where  $\text{Cor}(\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{a \times b}$  is the cross-correlation matrix between random vectors  $\mathbf{X}$  and  $\mathbf{Y}$ , and  $\odot$  represents the element-wise product of matrices (i.e., the Hadamard product).

*Proof.* Since we are concerned with covariances, we can assume by mean centering that  $\mathbf{E}\mathbf{X} = 0 \in \mathbb{R}^a$  and  $\mathbf{E}\mathbf{Y} = \mathbf{E}\epsilon = 0 \in \mathbb{R}^b$ . Note that  $D^t e_a = d e_b$  and  $U_a = e_a e_a^t$ , where  $e_r \doteq (1, \dots, 1)^t \in \mathbb{R}^r$  for  $r \in \{a, b\}$ . Hence using independence of  $\mathbf{X}$  and  $\epsilon$ :

$$\begin{aligned} \text{Cov}(\mathbf{Y}) &= \mathbf{E}(\mathbf{Y}\mathbf{Y}^t) = D^t \mathbf{E}(\mathbf{X}\mathbf{X}^t) D + \mathbf{E}(\epsilon \epsilon^t) \\ &= D^t \text{Cov}(\mathbf{X}) D + \text{Cov}(\epsilon) = D^t (\rho e_a e_a^t + (1 - \rho) I_a) D + \sigma^2 I_b \\ &= \rho (D^t e_a)^t (D^t e_a) + (1 - \rho) D^t D + \sigma^2 I_b \\ &= \rho d^2 e_b e_b^t + (1 - \rho) D^t D + \sigma^2 I_b \end{aligned}$$

Since all the diagonal entries of  $D^t D$  have the value  $d$ ,

$$\text{diag}[\text{Cov}(\mathbf{Y})] = (\rho d^2 + (1 - \rho) d + \sigma^2) I_b = (d\delta + \sigma^2) I_b = \left(\frac{\delta}{\eta}\right)^2 I_b \quad (4.5)$$

where for any square matrix  $A$ ,  $\text{diag}[A]$  denotes the diagonal matrix obtained from  $A$  by setting all the off-diagonal entries of  $A$  to 0.

We can similarly calculate the cross-covariance between  $\mathbf{X}$  and  $\mathbf{Y}$

$$\begin{aligned} \text{Cov}(\mathbf{X}, \mathbf{Y}) &= \mathbf{E}(\mathbf{X}\mathbf{Y}^t) = \mathbf{E}(\mathbf{X}\mathbf{X}^t) D = (\rho e_a e_a^t + (1 - \rho) I_a) D \\ &= \rho d e_a e_b^t + (1 - \rho) D, \end{aligned} \quad (4.6)$$

and also finally the cross-correlation between  $\mathbf{X}$  and  $\mathbf{Y}$  using (4.6), (4.5) and  $\text{diag}[\text{Cov}(\mathbf{X})] = I_a$ :

$$\begin{aligned} \text{Cor}(\mathbf{X}, \mathbf{Y}) &= \text{diag}[\text{Cov}(\mathbf{X})]^{-\frac{1}{2}} \text{Cov}(\mathbf{X}, \mathbf{Y}) \text{diag}[\text{Cov}(\mathbf{Y})]^{-\frac{1}{2}} \\ &= \frac{\eta}{\delta} (\rho d e_a e_b^t + (1 - \rho) D) = \frac{\eta}{\delta} (\rho d \bar{D} + (1 - \rho + \rho d) D) \\ &= \eta D + \eta \rho d \delta^{-1} \bar{D}. \end{aligned}$$

where  $\bar{D} \doteq 1 - D$ . In particular this shows (4.4). □

## 4.2 Real data study: bimodules for eQTL analysis

Now we describe the application of bimodules to the problem of expression quantitative trait loci (eQTL) analysis discussed in Section 1.1. The NIH funded GTEx Project has collected and created a large eQTL database containing genotype and expression data from postmortem tissues of human donors. A unique feature of this database is that it contains expression data from many tissues. We applied BSP, CONDOR and standard eQTL-analysis to  $p = 556,304$  SNPs and  $q = 26,054$  thyroid expression measurements from  $n = 574$  individuals. A detailed account of data acquisition, preprocessing, and covariate correction can be found in Section 4.3.2.

The 556K SNPs considered were a representative subset chosen from 4.9 million (directly observed and imputed) autosomal SNPs with minor allele frequency greater than 0.1. Using a representative set decreased computation time and reduced the multiple testing burden in each iteration of BSP. As SNPs exhibit local correlation due to linkage disequilibrium (LD), the selection process should not reduce the statistical power of BSP. We used an LD pruning software *SNPRelate* (135) to select the representative subset of SNPs (see Section 4.3.2 for details).

### 4.2.1 Results of BSP

We applied BSP to the thyroid eQTL data with false discovery parameter  $\alpha = 0.03$  selected to keep the edge-error under 0.05 (see Section 4.3.3). The search was initialized from singleton sets of all genes and half of the available SNPs, chosen at random. Thus the search procedure in Section 3.3.3 was run  $p/2 + q \sim 304\text{K}$  times. BSP took 4.7 hours to run on a computer with a 20-core 2.4 GHz processor (machine details provided in Section 4.3.4). The search identified 3744 unique bimodules with p-values below the significance threshold of  $\frac{\alpha}{pq} = 3.45 \times 10^{-12}$  (see Section 3.3.3). The majority (277K) of the searches terminated in the empty set after the first step; of the remaining 27K searches, the great majority identified a non-empty fixed point within 20 steps. Only 20 searches cycled and did not terminate in a fixed point. Among the searches taking more than one iteration, 94% terminated by the fifth step. Among searches that found a non-empty fixed point, 92.3% of the fixed points contained the seed singleton set of the search.

Among the unique bimodules discovered by BSP, some bimodules were similar to others; hence the effective number (3.4.4) of bimodules was 3304, slightly smaller than the number of unique

bimodules. We then applied the filtering procedure described in Section 3.4.4 to select a sub-collection of 3304 bimodules that were roughly disjoint. The selected bimodules had SNP sets ranging in size from 1 to 1000, and gene sets ranging in size from 1 to 100 (Figure 4.3); the median size of the gene and SNP sets was 1 and 7, respectively.

If required, BSP can be run in a faster (less exhaustive) or slower (more exhaustive) fashion by selecting a smaller or larger fraction of SNPs from which to initialize the search procedure. The effective number of discovered bimodules was only slightly smaller (3258) when initializing with 10% of the SNPs.

## 4.2.2 Running other methods

Standard eQTL analysis was performed by applying Matrix-eQTL (114) twice to the data, first to perform a *cis*-eQTL analysis within a distance of 1MB and next to perform a *trans*-eQTL analysis. In each case, SNP-gene pairs with BH (11)  $q$ -value less than 0.05 were identified as significant. Matrix-eQTL identified 186K *cis*-eQTLs and 73K *trans*-eQTLs.

To obtain CONDOR bimodules (107), we applied Matrix-eQTL to identify both *cis*- and *trans*-eQTLs with BH  $q$ -value under the threshold .2, chosen as in Fagny et al. (47). The resulting gene-SNP bipartite graph formed by these eQTLs was passed through CONDOR’s bipartite community detection pipeline (107), which partitioned the nodes of largest connected component of this graph into 6 bimodules.

We also applied the sCCA method of (130) using the permutation based parameter selection procedure (129) on the covariate-corrected genotype and expression matrices to identify 50 bimodules. The identified bimodules were large, containing roughly 100K SNPs and 4K-8K genes (see Figure 4.3), making them difficult to analyze and interpret. The identified bimodules also exhibited moderate overlap: the effective number was 25. As such, we excluded the sCCA bimodules from subsequent comparisons. Analysis of sCCA on the simulated data (Section 4.1.3) suggests that the method may be able to recover smaller bimodules with a more tailored choice of its parameters. While this is an interesting topic for future research, it is beyond the scope of the present paper.

### 4.2.3 Quantitative validation

In this subsection, we apply several objective measures to validate and understand the bimodules found by BSP and CONDOR.

#### 4.2.3.1 Permuted data

In order to assess the propensity of each method to detect spurious bimodules, we applied BSP and CONDOR to five data sets obtained by jointly permuting the sample labels for the expression measurements and most covariates (all except the five genotype PCs), while keeping the labels for genotype measurements and genotype covariates unchanged. Each data set obtained in this way is a realization of the permutation null defined in Definition 3. BSP found very few (5-12) bimodules in the permuted datasets compared to the real data (3304). CONDOR found no bimodules in any of the permuted datasets.

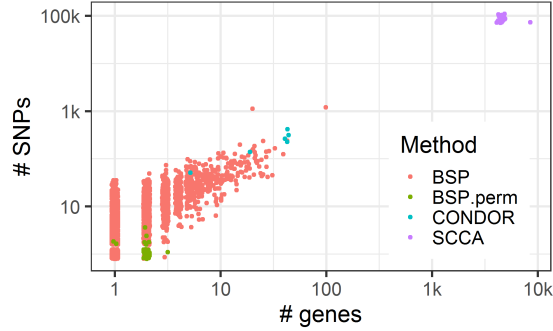
#### 4.2.3.2 Bimodule sizes

Most (89%) bimodules found by BSP have fewer than 4 genes and 50 SNPs, but BSP also identified moderately sized bimodules having 10-100 genes and 30-1000 SNPs (see Figure 4.3). The bimodules found by CONDOR were moderately sized, with 10-100 genes and several hundred SNPs, except for a smaller bimodule with 5 genes and 43 SNPs. On the permuted data, most bimodules found by BSP have fewer than 2 genes and 2 SNPs.

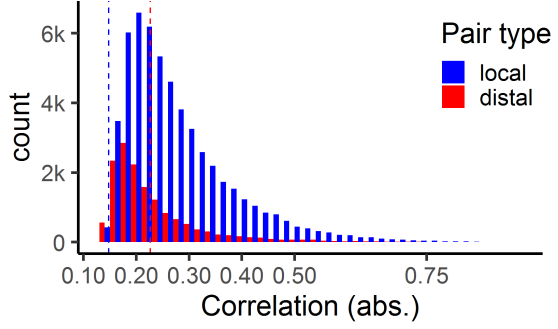
As a one dimensional measure, we define the *geometric size* of a bimodule  $(A, B)$  to be  $\sqrt{|A||B|}$ , the geometric mean of its gene and SNP counts, or equivalently, the square root of the number of gene-SNP pairs in the bimodule.

#### 4.2.3.3 Connectivity threshold and network sparsity

Stable bimodules capture aggregate association between groups of SNPs and genes, however it is unclear how to recover individual SNP-gene associations within these bimodules. Motivated by the network perspective, in Section 3.3.4 we proposed evaluating for each bimodule  $(A, B)$ , the connectivity threshold (3.12) and the corresponding network of essential edges (3.13) between  $A$  and  $B$ . To understand the structure of the network of essential edges, we further calculated the



**Figure 4.3:** The sizes of bimodules detected by BSP, CONDOR and sCCA, and sizes of bimodules detected by BSP under the 5 permuted datasets.



**Figure 4.4:** Correlations corresponding to SNP-gene pairs that appear as essential-edges (Section 3.3.4) in one or more BSP bimodules with geometric size above 10. Local pairs to the left of the blue line (*cis*-analysis threshold) and distal pairs to the left red line (*trans*-analysis threshold) show importance at the network level but were not discovered by standard eQTL analysis.

*tree-multiplicity*

$$\text{TreeMul}(A, B) \doteq \frac{|\text{essential-edges}(A, B)|}{|A| + |B| - 1}, \quad (4.7)$$

which measures the number of essential edges relative to the number of edges in a tree on the same node set.  $\text{TreeMul}(A, B)$  is never less than 1, and takes the value 1 exactly when the essential edges form a tree.

For bimodules found by BSP, the connectivity thresholds ranged from 0.14 to 0.59 and tree-multiplicities ranged from 1 to 10; the smaller values of the former and larger values of latter were associated with bimodules of larger geometric size (Figure 4.7). Smaller bimodules had large connectivity thresholds and a tree-like essential edge network; in other words, such bimodules were connected under a small number of strong and local (see Section 4.2.4.2) SNP-gene associations. On the other hand larger bimodules had lower connectivity thresholds, meaning that we had to include weaker and often distal (see Section 4.2.4.2) SNP-gene associations to connect such bimodules. After including the weaker SNP-gene edges, although the association network for large bimodules had tree-multiplicity around 10 (Figure 4.7), these networks were still sparsely connected compared to the complete bipartite graph on the same nodes.



#### 4.2.4 Biological Validation

In order to assess potential biological utility of bimodules found by BSP, we compared the SNP-gene pairs in bimodules to those found by standard *cis*- and *trans*-eQTL analysis, studied the locations of the SNPs, and examined the gene sets for enrichment of known functional categories.

##### 4.2.4.1 Comparison with standard eQTL analysis

As described earlier, the bimodules produced by CONDOR are derived directly from SNP-gene pairs identified by *cis*- and *trans*-eQTL analysis. Table 4.1 compares these eQTL pairs with those found in bimodules identified by BSP. Recall that *cis*-eQTL analysis considers only local SNP-gene pairs (improving detection power by reducing multiple testing), while *trans*-eQTL analysis and BSP do not use any information about locations of the SNPs and genes. We find that half of the pairs identified by *cis*-eQTL analysis and most of the pairs identified by *trans*-eQTL analysis appear in at least one bimodule.

Bimodules capture sub-networks of SNP-gene associations rather than individual eQTLs, and as such individual SNP-gene pairs in a bimodule need not be eQTLs. In fact, the results of Section 4.2.3.3 suggest that the association networks underlying large bimodules may be sparse. Define a bimodule  $(A, B)$  to be connected by a set of eQTLs if the bipartite graph with vertex set  $A \cup B$  and edges corresponding to the eQTLs is connected. As shown in Table 4.1, a significant fraction of BSP bimodules are not connected by either *cis*- or *trans*-eQTLs. The discovery of such bimodules suggests that the sub-networks identified by BSP cannot be found by standard eQTL analysis, and that these sub-networks can provide new insights and hypotheses for further study.

To identify potentially new eQTLs using BSP, we examine bimodule connectivity under the combined set of *cis*- and *trans*-eQTLs. All of the bimodules with one SNP or one gene are connected by the combined set of eQTLs (Section 4.3.6), and therefore all edges in these bimodules are discovered by standard analyses. On the other hand, 224 out of the 358 bimodules with geometric size larger than 10 were not connected by the combined set of eQTLs. In Figure 4.4, we plot the correlations corresponding to SNP-gene pairs that appear as essential-edges (Section 3.3.4) in one or more bimodules with geometric size above 10, along with the correlation thresholds for *cis*-eQTL (blue line) and *trans*-eQTL (red line) analysis. Around 300 local edges (i.e. the SNP is located

Analysis type	% eQTLs found among bimodules	% bimodules connected by eQTLs
<i>trans</i> -eQTL analysis	84%	70%
<i>cis</i> -eQTL analysis	51%	88%

**Table 4.1:** Comparison of BSP and standard eQTL analysis. A gene-SNP pair is said to be found among a collection bimodules if the gene and SNP are both part of some common bimodule. On the other hand, we say that a bimodule is connected by a collection of eQTLs if under the gene-SNP pairs from the collection, the bimodule forms a connected graph.

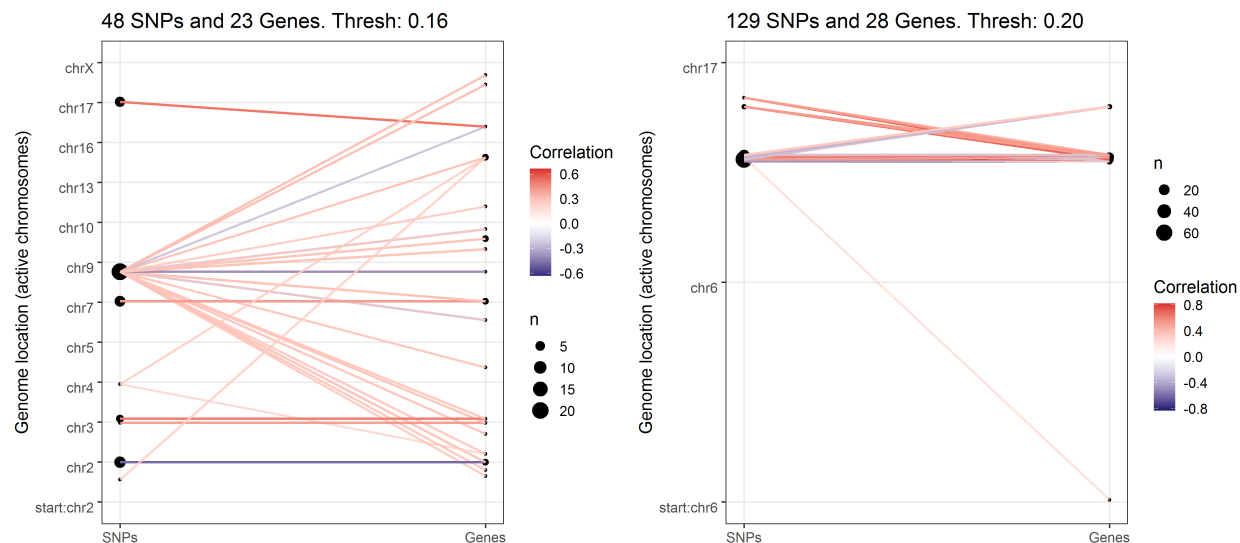
within 1MB of the gene transcription start site) and 8.8K distal edges do not meet the correlation thresholds for *cis*- and *trans*-eQTL analysis, respectively, but show evidence of importance at the network level, and may be worthy of further study.

#### 4.2.4.2 Genomic locations

We studied the chromosomal location and proximity of SNPs and genes from bimodules found by BSP and CONDOR. While CONDOR uses genomic locations as part of the *cis*-eQTL analysis in its first stage, BSP does not make use of location information. Genetic control of expression is often enriched in a region local to the gene (37). All CONDOR bimodules, and almost all (99.3%) BSP bimodules, have at least one local SNP-gene pair (the SNP is located within 1MB of the gene transcription start site). In 93.5% of the smaller BSP bimodules (geometric size 10 or smaller) and 54.8% of the medium to large BSP bimodules (geometric size above 10) each gene and each SNP had a local counterpart SNP or gene within the bimodule.

For each bimodule, we examined the chromosomal locations of its SNPs and genes. All SNPs and many of the genes from the six CONDOR bimodules were located on Chromosome 6; two CONDOR bimodules also had genes located on Chromosome 8 and Chromosome 9. The SNPs and genes from the BSP bimodules were distributed across all 23 chromosomes: 170 of the 2947 small bimodules spanned 2 to 5 chromosomes and 152 of the 358 medium to large bimodules spanned 2 to 11 chromosomes; however the remaining bimodules were localized to a chromosome each.

Figure 4.5 illustrates the genomic locations of two bimodules found by BSP, with SNP location on the left and gene location on the right (only active chromosomes are shown). In addition, the figure illustrates the essential edges (Section 3.3.4) of each bimodule. The resulting bipartite graph provides insight into the underlying associations between SNPs and genes that constitute the bimodule. See Section 4.3.7 for more such illustrations.



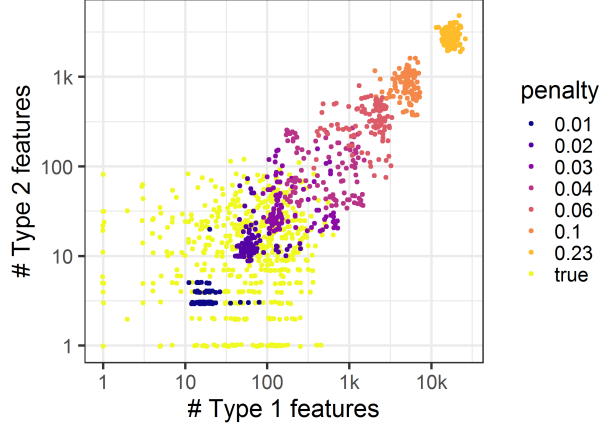
**Figure 4.5:** The gene-SNP association network for two BSP bimodules mapped onto the genome. The network of *essential edges* was formed by thresholding the cross-correlation matrix for the bimodule at the connectivity threshold (Section 3.3.4).

#### 4.2.4.3 Gene Ontology enrichment for bimodules

The Gene Ontology (GO) database contains a curated collection of gene sets that are known to be associated with different biological functions (c.f. 50, 20, 111). The topGO (2, 3) package assesses whether sets in the GO database are enriched for a given gene set using Fisher’s test. For each of the 145 BSP bimodules having a gene set  $B$  with 8 or more elements, we used topGO to assess the enrichment of  $B$  in 6463 GO gene sets of size more than 10, representing biological processes; however these significant sets were not apparently related to thyroid-specific function. We retained results with significant BH  $q$ -values ( $\alpha = .05$ ). Of the 145 gene sets considered, 18 had significant overlap with one or more biological process. Repeating with randomly chosen gene sets of the same size yielded no results. The significant GO terms for BSP and CONDOR can be found in Appendix A.

### 4.3 Details of data analysis

This final section of the chapter contains various details about the data analysis conducted earlier. In particular, Section 4.3.1, presented next, contains details about the simulation study in



**Figure 4.6:** The sizes of sCCA bimodules for various values of the penalty parameter  $\lambda$  along with sizes of the true bimodules.

Section 4.1, while the remaining sections present details about the real data application in Section 4.2.

#### 4.3.1 Results from sCCA on the simulation study

We ran sCCA on the simulated dataset to search for 100 canonical covariates for a range of values of the penalty parameter  $\lambda$ . The sizes of the bimodules for various values of  $\lambda$  can be seen in Figure 4.6. For  $\lambda \in \{0.01, 0.02, 0.03, 0.04, .233\}$ , the first two columns of the following table show the number of true bimodules (TB) that overlapped with each detected bimodule (DB) and the edge-error of each DB averaged over all DBs. The last column shows the top 5 (or bottom 95) percentile *recall* among the true bimodules.

$\lambda$	# TBs that overlap with each DB	edge-error	recall of TB (95%-tile)
.01	.97	.09	.02
.02	.96	.19	.23
.03	1.97	.48	.62
.04	6.47	.65	.95
.23	281	.89	1

The parameter value  $\lambda = 0.01$  has small edge-error, but poor recall. The recall improves on increasing  $\lambda$ , but the edge-error degrades.

### 4.3.2 Data acquisition and preprocessing for the GTEx eQTL data

We obtained genotype and thyroid expression data for 574 individuals from the dbGap website (accession number: phs000424.v8.p1). We directly used the filtered and normalized gene expression data and covariates provided for eQTL analysis but filtered the SNPs in the genotype data using the LD pruning software *SNPRelate* (135). The software retained 556K autosomal SNPs with minor allele frequency above 0.1 such that all pairs of SNPs within each 500KB window of the genome had squared correlation under  $(0.7)^2$ . The latter threshold was chosen to balance the number of retained SNPs and information loss.

There were 68 covariates provided for the Thyroid tissue consisting of the top 5 genotype principle components; 60 PEER covariates, and 3 additional covariates for sequencing platform, sequencing protocol, and sex. We accounted for these covariates by the modification to BSP mentioned in Section 3.4.5.

### 4.3.3 Choice of false discovery parameter to BSP

We chose the false discover parameter  $\alpha$  for BSP from the grid  $\{0.01, 0.02, 0.03, 0.04, 0.05\}$  by finding the largest  $\alpha$  that kept the average edge-error estimates based on  $N = 5$  half-permutations under 0.05 (see Section 3.4.3). However our error estimates were variable as we obtained  $\alpha = 0.05$  in one instance and  $\alpha = 0.03$  in another. We conservatively chose  $\alpha = 0.03$ .

### 4.3.4 Details of our hardware and software stack

The various methods used in this analysis were run on a dedicated computer that had Intel (R) Xeon (R) E5-2640 CPU with 20 parallel cores at 2.50 Hz base frequency, and a 512 GB random access memory along with L1, L2 and L3 caches of sizes 1.3, 5 and 50 MB respectively. The computer ran Windows server 2012 R2 operating system and we used the Microsoft R Open 3.5.3 software to perform most of our analysis, since it has multi-core implementations of linear algebra routines.

### 4.3.5 Bimodule connectivity thresholds and network sparsity

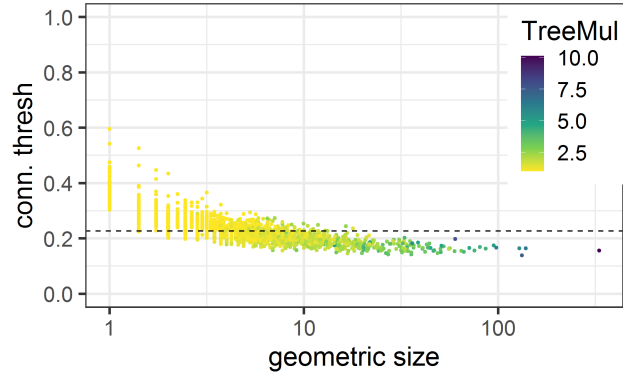
Figure 4.7 shows the two network statistics for bimodules found by BSP – connectivity threshold and tree-multiplicity. All bimodules have tree multiplicity under 10. This shows that the association network for large bimodules, particularly having low connectivity-thresholds, is sparse.

### 4.3.6 Connectivity of BSP bimodules under standard eQTL analysis

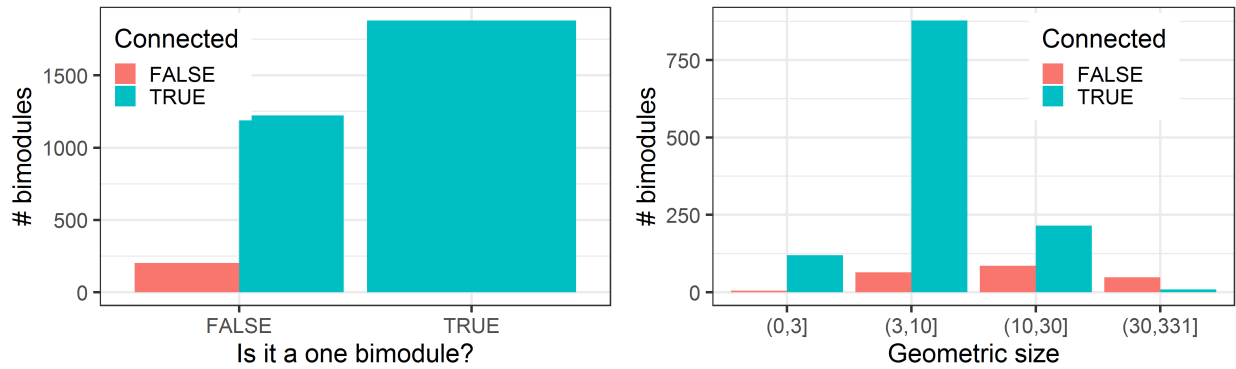
Here we examine which bimodules are connected under the combined edges from *cis*-eQTL and *trans*-eQTL analysis, based on geometric size of the bimodule. Figure 4.8 (left) shows that all the bimodules that have either one gene or one SNP are connected. Hence, these bimodules could have been recovered using standard eQTL analysis. On the other hand if we restrict to bimodules with two or more genes and SNPs, we see that (Figure 4.8; right) the fraction of connected bimodules tends to decrease as the geometric size of the bimodules increases.

### 4.3.7 Genomic plots of more BSP bimodules

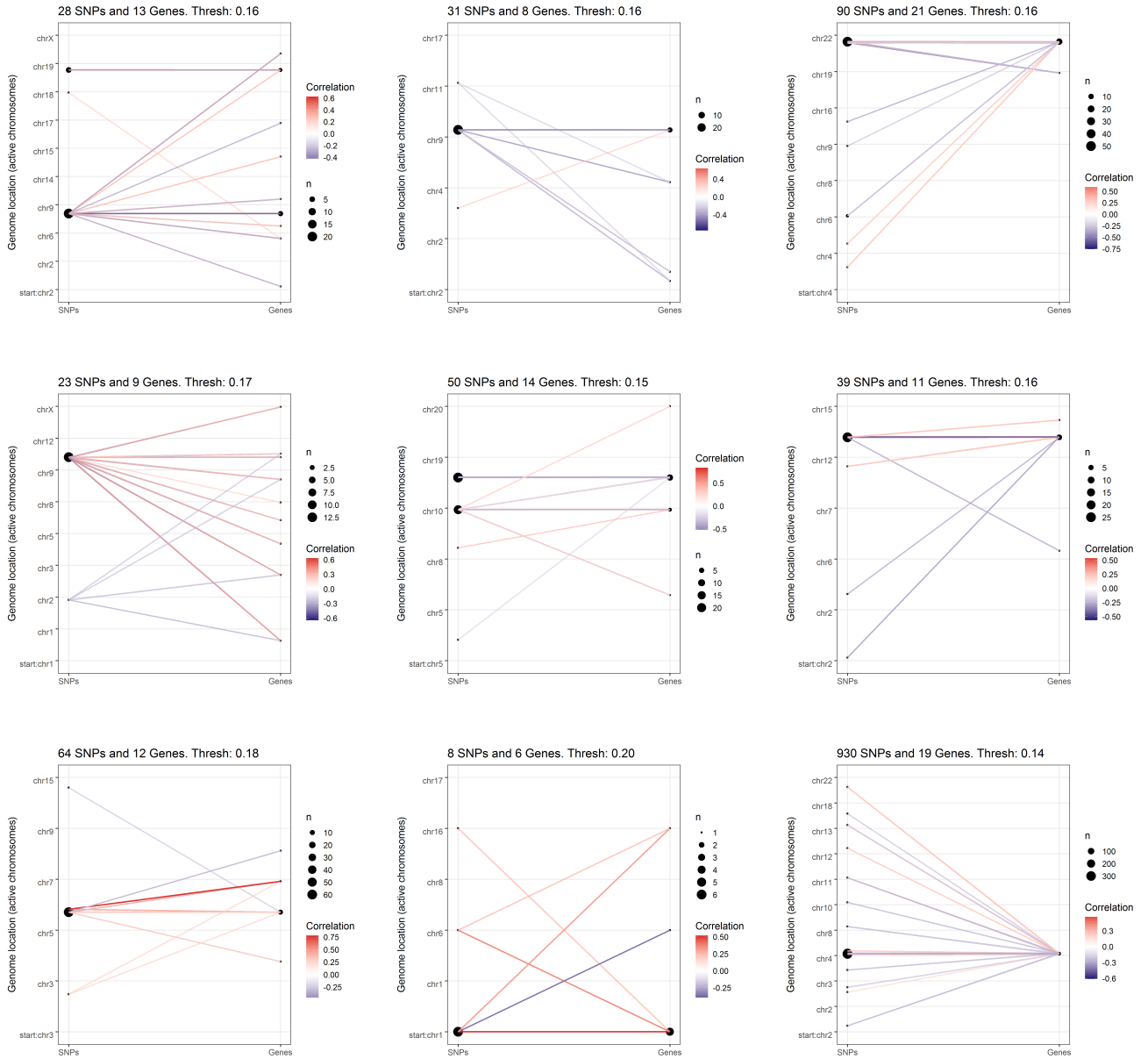
See the plots in Figure 4.9.



**Figure 4.7:** Connectivity-threshold and tree-multiplicity for BSP bimodules compared to their geometric size. The horizontal dotted line represents the threshold obtained from standard trans-analysis.



**Figure 4.8:** Connectivity of BSP bimodules under combined edges from *cis*-eQTL and *trans*-eQTL analysis. Left: the number of bimodules that are connected and are *one bimodules*, i.e. have one gene or one SNP. Right: Among bimodules having two or more genes and SNPs (i.e. are not one bimodules), the connectivity and geometric size of the bimodules.



**Figure 4.9:** Out of 31 BSP bimodules that had genes on 3 or more chromosomes and SNPs on 2 or more chromosomes, we selected 9 bimodules that looked interesting. The bipartite graph for each bimodule is formed out of the essential edges (Section 4.2.3.3).



## CHAPTER 5

### Literature review for distributed load balancing

In this chapter we will review literature related to the power-of-choice and distributed load balancing. In Section 5.1, we start with the simple setup of the balls and bins problem. These are typically static problems studied in computer science where the system load is analyzed after all the balls have been placed into bins.

Next, in Section 5.2, we review literature on continuous time versions of the balls and bins problem in the setup of queueing theory, where balls (now jobs) arrive and depart (i.e. they are served) from the system. These problems have recently garnered interest from both theoretical and applied researchers. From an applied perspective, they can be used to study the performance of a variety of load-balancing algorithm in the setup of modern large-scale data centers. From a theoretical perspective, the analysis of such systems draws from a variety of fields of probability theory like stochastic analysis, interacting particle systems, Stein’s method, and coupling arguments. Such techniques, for instance, allow various aspects of the load balancing systems to be analytically approximated when the system size is large.

We now provide a survey of these results.

#### 5.1 The balls and bins problem

Consider the following balls and bins problem briefly mentioned in Chapter 1: suppose  $n$  balls are to be sequentially placed into  $n$  bins with the objective of limiting the maximum bin size after all the balls have been placed (where we define the size of a bin to be the number of balls that are currently placed in that bin). In the absence of a central dispatcher, each incoming ball must communicate with a bin to determine its current size. Hence, in order to reduce the communication overhead, suppose that the arriving ball chooses to inspect the sizes of only  $d \in [n]$  randomly chosen bins. If  $d \in [n]$  is the same for each ball and the balls choose to be placed into the smallest bin

among the  $d$  inspected bins, this is called the power-of- $d$  assignment scheme. This is a ‘distributed’ scheme because each ball makes its own decisions.

For the power-of- $d$  scheme with a fixed  $d \in [n]$  and large  $n$ , with high probability, the maximum bin size after all the balls has been placed is  $(1 + o(1)) \frac{\log n}{\log \log n}$  for  $d = 1$  (see e.g. (110)) and  $\log_d \log n + \Theta(1)$  for any fixed  $d \geq 2$  (6). This substantial decrease in the maximum bin size on moving from purely random assignment ( $d = 1$ ) to random assignment with limited choice (fixed  $d \geq 2$ ) is often called as the ‘the power of choice’. This phenomenon has numerous applications to areas like large scale load balancing, hashing, collision protocols, network exploration, and distributed learning (112, 127, 5, 51, 35).

Similar results have also been studied when, in a more general setup,  $m$  balls have to be placed into  $n$  bins. The standard proof techniques extend to the case  $m = O(n)$ , but the so-called ‘heavily loaded case’ where  $m \gg n$ , has been more challenging. Nevertheless, when  $m \geq n \log n$ , it was shown that maximum load is  $m/n + \Theta(\sqrt{m/n \log n})$  when  $d = 1$  (110) and  $m/n + \log_d \log(n) + \Theta(1)$  when  $d \geq 2$  (13, 118), with high probability. Hence power-of- $d$  scheme for  $d \geq 2$  satisfies the interesting property that the difference between the maximum bin size and the average bin size for the balls and bins system is, with high probability, independent the total number of balls  $m$  in the system.

Various further extensions of the above balls and bins problem has been studied. This includes (see (14, 100) and references therein) models with parallel and batch allocation of balls, infinite processes where balls are continuously added and removed, and models where the selection of bins is constrained in terms of a graph (108, 104, 64). The power-of- $d$  scheme described earlier is a non-adaptive strategy since the number of bins inspected by each ball is always  $d$ . Adaptive strategies, where the number of choices for a ball may increase based on the currently inspected load, can have better performance when compared to the power-of- $d$  scheme (40). In fact, some (15, 69) adaptive strategies, can achieve a constant maximum bin load using only an average of  $O(1)$  bin choices per ball. However, compared to the power-of- $d$  scheme, these strategies need additional assumptions like multiple rounds of communication, or additional information like the current number of balls in the system.

### 5.1.1 The $(k, d)$ balls and bins problem

In this dissertation, we focus on the asymptotic performance of power-of- $d$  load balancing scheme when  $d = d_n$  is allowed to increase with  $n$ . In fact, in Chapter 6, for any integer  $1 \leq k < d \leq n$ , we consider a generalization of the traditional balls and bins problem described above, called  $(k, d)$  balls and bins problem. In this problem,  $m = nk$  balls are sequentially placed into  $n$  bins in batches of size  $k$ . For this, at each time  $t \in [n]$ ,  $d$  bins are chosen at random, and an incoming batch of  $k$  balls are placed into the  $k$ -smallest of the  $d$  chosen bins. Hence the case  $k = 1$  is just the standard balls and bins problem with the power-of- $d$  scheme. The general case of  $k > 1$  is motivated by problems like scheduling jobs with  $k$  parallel tasks in cluster environments (97), and cloud storage applications where parts of each file may be coded and stored on multiple servers so that any subset of  $k$  distinct parts be used to recover the entire file (70). A comprehensive analysis of the maximum for the  $(k, d)$  balls and bins problem has appeared in (100), but as a warm up, in Chapter 6 we present our calculations for the same problem, noting throughout that  $(k, d) = (k(n), d(n))$  may now depend on  $n$ .

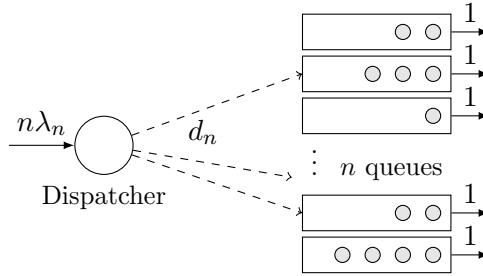
## 5.2 The Supermarket model

Consider a processing system with  $n$  parallel queues in which each queue's jobs are processed by the associated server at rate 1. Jobs arrive at rate  $n\lambda_n$  and join the shortest queue amongst  $d_n$  randomly selected queues (without replacement), with  $d_n \in [n]$  (see Figure 5.1). The interarrival times and service times are mutually independent exponential random variables. This queuing system with the above described 'join-the-shortest-queue amongst chosen queues' discipline is often denoted as  $JSQ(d_n)$  and frequently referred to as the Supermarket model (cf. (52, 79, 78, 75, 82, 91) and references therein). Note that when  $d_n = n$  the above description corresponds to a policy where an incoming job joins the shortest of all queues in the system (see e.g. (45)). The case  $d_n = 1$  is the other extreme corresponding to incoming jobs joining a randomly chosen queue in which case the system is equivalent to one with  $n$  independent  $M/M/1$  queues with arrival rate  $\lambda_n$  and service rate 1. The case  $d_n = d$  where  $d > 1$  is a fixed positive integer is sometimes also referred to as the power-of- $d$  scheme. The analysis of  $JSQ(d_n)$  schemes has been a focus of much recent research motivated by problems from large scale service centers, cloud computing platforms, and data storage and

retrieval systems (see e.g. (29, 122, 80, 96, 92, 55, 4, 22)). The influential works of Mitzenmacher (89, 112) and Vvedenskaya et al. (124) showed by considering a fluid scaling that increasing  $d$  from 1 to 2 leads to significant improvement in performance in terms of steady-state queue length distributions in that the tails of the asymptotic steady-state distributions decay exponentially when  $d = 1$  and super-exponentially when  $d = 2$ . Additionally, concentration and rapid mixing results from Luczak and McDiarmid (78) showed that when  $\lambda_n = \lambda < 1$  is fixed, with high probability, the maximum queue length in steady state is  $(1 + o(1)) \frac{\log n}{\log(1/\lambda)}$  for  $d = 1$  and  $\frac{\log \log n}{\log d} + O(1)$  for any fixed  $d \geq 2$ , as  $n \rightarrow \infty$ . Strong approximations and limit theorems under a diffusion scaling for the  $JSQ(d)$  system, with a fixed  $d$ , can be found in (53, 44, 26, 79). Although  $JSQ(d)$  for a fixed  $d \geq 2$  leads to significant improvements over  $JSQ(1)$ , as observed in (48, 49), no fixed value of  $d$  provides the optimal waiting time properties of the join-the-shortest-queue system (i.e.  $JSQ(n)$ ) particularly in the heavy traffic regime when  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta > 0$ . Although there are other load balancing algorithms (see the survey (123)) similar to Section 5.1 that can use additional information to overcome limitations of  $JSQ(d)$ , we will focus on the behavior of the  $JSQ(d_n)$  system in heavy traffic (i.e.  $\lambda_n \uparrow 1$ ) as  $d_n$  increases with the system size  $n$ . Such an asymptotic study is carried out in Chapter 7, but here we review the existing work on  $JSQ(d_n)$  systems as  $d_n \rightarrow \infty$ . Table 1.1 in Chapter 1 provides a comparison of the regimes covered by existing work.

### 5.2.1 Universal law of large numbers as $d_n \rightarrow \infty$

The paper Mukherjee et al. (91) studied the law of large numbers (LLN) behavior of a  $JSQ(d_n)$  system, under a standard scaling, when  $d_n \rightarrow \infty$ . The precise result of (91) is as follows. For  $i \in \mathbb{N}_0 \doteq \{0, 1, 2, \dots\}$  and  $t \in [0, \infty)$ , let  $G_{n,i}(t)$  denote the fraction of queues with at least



**Figure 5.1:** Illustration of the Supermarket model with parameters  $(n, d_n, \lambda_n)$ . There are  $n$  queues, each having a server processing jobs at rate 1. Jobs arrive to the system at rate  $n\lambda_n$  and independently sample  $d_n \in \{1, \dots, n\}$  of the queues at random, and enter the smallest of the  $d_n$  selected queues.

$i$  customers at time  $t$  in the  $n$ -th system. Note that  $G_{n,0}(t) = 1$  for all  $t \geq 0$ . We will call  $\mathbf{G}_n(t) \doteq (G_{n,i}(t))_{i \geq 1}$  the state occupancy process. This process has sample paths in the space of summable nonnegative sequences. More precisely, for  $p \geq 1$ , let  $l_p$  be the space of real sequences  $\mathbf{x} := (x_1, x_2, \dots)$  such that  $\|\mathbf{x}\|_p \doteq (\sum_{i=1}^{\infty} |x_i|^p)^{1/p} < \infty$ . Let

$$l_1^\downarrow \doteq \{\mathbf{x} \in l_1 : x_i \geq x_{i+1} \text{ and } x_i \in [0, 1] \text{ for all } i \in \mathbb{N}\} \quad (5.1)$$

be the space of non-increasing sequences in  $l_1$  with values in  $[0, 1]$ , equipped with the topology generated by  $\|\cdot\|_1$ . Note that  $l_1^\downarrow$  is a closed subset of  $l_1$  and hence is a Polish space. Then, whenever  $\|\mathbf{G}_n(0)\|_1 < \infty$  a.s., it can be shown that  $\{\mathbf{G}_n(t) : t \geq 0\}$  is a stochastic process with sample paths in  $\mathbb{D}([0, \infty) : l_1^\downarrow)$  (the space of right continuous functions with left limits from  $[0, \infty)$  to  $l_1^\downarrow$  equipped with the usual Skorohod topology); see Section 7.3. The paper (91) shows the following two facts under the assumption that  $\mathbf{G}_n(0)$  converges in probability to some  $\mathbf{r} \in l_1^\downarrow$ :

- (a) When  $d_n = n$  and  $\lambda_n \rightarrow \lambda \in (0, \infty)$ ,  $\mathbf{G}_n$  is a tight sequence in  $\mathbb{D}([0, \infty) : l_1^\downarrow)$  and every weak limit point satisfies a certain set of “fluid limit equations” (see Theorem 5 in (91), and equations (7.5)-(7.6) in Chapter 7);
- (b) When  $d_n$  is an arbitrary sequence growing to  $\infty$  and  $\lambda_n \rightarrow \lambda \in (0, 1)$ , then the statements in (a) hold once more for  $\mathbf{G}_n$ .

Our work in Chapter 7 begins by revisiting the above LLN results from (91). In Theorem 8 in Chapter 7, we show that, when  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{r}$ , for arbitrary sequences  $d_n \rightarrow \infty$  and  $\lambda_n \rightarrow \lambda \in (0, \infty)$ ,  $\mathbf{G}_n$  converges in probability in  $\mathbb{D}([0, \infty) : l_1^\downarrow)$  to a continuous trajectory  $g$  in  $l_1^\downarrow$  that is characterized as the *unique* solution of an infinite system of constrained ordinary differential equations (ODE) (see (7.3) in Proposition 7 in Chapter 7). Using standard properties of the Skorohod map we observe in Remark 7.2.2 that a continuous trajectory in  $l_1^\downarrow$  solves the fluid limit equations of (91) if and only if it solves (7.3). This together with Proposition 7 proves that the fluid limit equations in (91) in fact have a unique solution. In this manner we complete and strengthen the result from (91). Our proof of the LLN result is quite different from the arguments in (91). The latter are based on sophisticated ideas of separation of time scales and weak convergence of measure valued processes from (62) to handle the convergence for  $d_n = n$ , and certain coupling techniques to treat

the general case when  $d_n < n$  and  $d_n \rightarrow \infty$ . In contrast, our approach is more direct and uses martingale estimates and well known characterization properties of solutions of Skorohod problems (see e.g. proof of Lemma 20).

### 5.2.2 Universal central limit theorem for $d_n \gg \sqrt{n} \log n$

Our main goal in Chapter 7 will be to study diffusion approximations for  $\mathbf{G}_n$  in the heavy traffic regime, namely when  $\lambda_n \rightarrow 1$ . In the case when  $d_n = n$  ( $JSQ(n)$  system), this problem has been studied in Eschenfeldt and Gamarnik (45). Their basic result is as follows. Suppose  $d_n = n$  and  $\sqrt{n}(\lambda_n - 1) \rightarrow \beta > 0$ . Consider the unit vector  $\mathbf{e}_1 = (1, 0, \dots)$  in  $l_2$ . Then under conditions on  $\mathbf{G}_n(0)$ ,  $\mathbf{Y}_n(\cdot) \doteq \sqrt{n}(\mathbf{G}_n(\cdot) - \mathbf{e}_1)$  converges in distribution in  $\mathbb{D}([0, \infty) : l_2)$  to a continuous stochastic process  $\mathbf{Y} = (Y_1, Y_2, \dots)$ , described in terms of a one dimensional Brownian motion, for which  $Y_i = 0$  for  $i > r$  for some  $r \in \mathbb{N}$  (which depends on the conditions assumed on  $\mathbf{G}_n(0)$ ). Specifically, when  $r = 2$ , the pair  $Y_1, Y_2$  is given as a two dimensional diffusion in the half space  $(-\infty, 0] \times \mathbb{R}$  with oblique reflection in the direction  $(-1, 1)^t$  at the boundary  $\{0\} \times \mathbb{R}$ . (For the form of the limit in the general case see Corollary 6, Chapter 7). In Mukherjee et al. (91) this result is extended to the case where  $d_n < n$  and  $d_n \gg \sqrt{n} \log n \rightarrow \infty$ . Under the same assumptions on the initial condition as in (45), it is shown in (91) that  $\mathbf{Y}_n$  converges to the same limit process as for the case  $d_n = n$ . The proof, as for the LLN result, proceeded by constructing a suitable coupling between a  $JSQ(d_n)$  and  $JSQ(n)$  system. The paper (91) also argued that when  $d_n \ll \sqrt{n} \log n$ , the process  $\mathbf{Y}_n$  cannot be tight and thus in this regime the above diffusion approximation cannot hold.

### 5.2.3 Steady state analysis as $d_n \rightarrow \infty$

Recent work has also considered the analysis of  $JSQ(d_n)$  systems in steady state as  $d_n \rightarrow \infty$ . For the case  $d_n = n$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta$ , Braverman (23) showed that the suitably scaled steady state distribution of  $JSQ(n)$  converges to the steady state of the limiting diffusion, mentioned above, from (45). Next, Banerjee and Mukherjee (7, 8) conducted a thorough analysis of the steady state distribution of the limiting diffusion, which could be used to obtain performance estimates for the  $JSQ(n)$  system in steady state.

The results mentioned so far have mainly applied to the  $JSQ(d_n)$  system in the Halfin-Whitt traffic regime where  $1 - \lambda_n = \Theta(1/\sqrt{n})$ . More generally, suppose we are interested in the traffic intensity  $\lambda_n = 1 - \gamma n^{-a}$  for some  $a \in (0, 1]$  and  $\gamma > 0$ . The choices  $a < 1/2$  and  $a > 1/2$  fall under the so-called sub- and super-Halfin-Whitt regimes, respectively. For the  $JSQ(d_n)$  system in this regime, Brightwell et al. (25, 24) prove that when  $d_n = n^b$  for some  $b \in (0, 1)$ , then with high probability, the system in steady state will be perfectly balanced and the size of the largest queue will be  $k = \lceil a/b \rceil$  (or one of  $\{k, k + 1\}$  if  $a/b = k$ ) if the condition  $2a < 1 + b(k - 1)$  holds. Liu and Ying (74) consider the case  $a \in (0, 1/2)$  and, when  $d_n \geq \gamma^{-1} n^a \log n$ , show that the  $JSQ(d_n)$  system with a slowly growing buffer size will have a vanishingly small fraction of jobs in steady state that have to wait for service.

## CHAPTER 6

### Maximum load in the balls and bins problem

In this chapter, as an instructive exercise, we estimate the maximum bin size in the balls and bins problem with the power-of- $d$  scheme and in the generalized  $(k, d)$  balls and bins problem, both of which were introduced in Chapter 5, Section 5.1. In our calculations, we will especially allow  $(k, d) = (k(n), d(n))$  to depend on  $n$  – the number of bins in the system. Although a tighter analysis of the maximum in the  $(k, d)$  balls and bins problem can be found in Park (100), the calculations in this chapter are an attempt to extend the concentration techniques from Luczak and McDiarmid (77) to the case where  $d = d(n)$  is allowed to depend on  $n$ .

#### 6.1 Model description and overview

Fix positive integers  $1 \leq k < d \leq n$ . The  $(k, d)$  balls and bins problem with  $n$  bins is the following random process. There are  $n$  bins, numbered from  $1, \dots, n$ , that are all initially (at time  $\tau = 0$ ) empty. Next, at each time  $\tau \in [n]$ ,  $d$  bins are chosen independently and uniformly at random (*without* replacement) out of these  $n$  bins, and one ball each is placed into the  $k$  smallest of these bins (the *size* of a bin is defined to be number of balls it currently contains). We assume that any ambiguity in the choice of the  $k$ -smallest bins is resolved by using a tie-breaking rule; for concreteness, we will assume that a smaller numbered bin is preferred over a larger numbered bin when both bins have the same size.

We would like to show that the maximum bin size at terminal time  $\tau = n$  is  $\frac{\log \log n}{\log(d-k+1)} + O(1)$  with high probability, under various regimes of  $(k, d) = (k(n), d(n))$  as  $n \rightarrow \infty$ . Although we will only be able to consider the cases  $(k(n), d(n)) = (1, O(\log n))$  or  $k(n) < d(n) = O(1)$  in this chapter, as shown in Park (100), this result holds more generally as long as  $d(n)/(d(n) - k(n)) = O(1)$ . Now we comment on proof idea in this chapter. First, in Section 6.2, we establish concentration results following (77) and derive its consequence that the maximum bin size at time  $n$  is one of



$\{m^*(n)-1, m^*(n), m^*(n)+1\}$  with high probability, for some constant  $m^*(n) \in \mathbb{N}$  determined by the mean behavior of the system. Hence to estimate  $m^*(n)$ , we study the mean behavior of the  $n$ -bin system via the “fluid limit” approximation in Section 6.3. Here, using quantitative error estimates, this mean behavior is approximated by the solution to a series of differential equations. However, the error analysis is complicated by the fact that the limiting differential equations (i.e. the fluid limit) depends on  $(k, d)$ , and we could not derive good explicit error bounds for the case  $k > 1$  (see (6.27) for an implicit bound). Finally, in Section 6.3.3, by analyzing the limiting differential equation when  $(k, d) = (1, O(\log n))$  or  $k < d = O(1)$ , we provide upper and lower estimates on its solution, which is then used to determine that  $m^*(n) = \frac{\log \log n}{\log(d-k+1)}$  upto bounded additive constants as  $n \rightarrow \infty$ .

### 6.1.1 Notation

Recall the random process associated with the  $(k, d)$  balls and bins problem with  $n$  bins, stated above. First in Section 6.2.1, the number of bins  $n$  will be fixed, and  $X_\tau \doteq (x_1^\tau, \dots, x_n^\tau) \in \mathbb{N}_0^n$  will be used to denote the state of the random process at time  $\tau \in [n]$ , where  $x_i^\tau \in \mathbb{N}_0$  represents the size of bin (i.e. the number of balls in the bin) labeled  $i \in [n]$  at time  $\tau$ . Define the associated filtration  $\mathcal{F}_\tau \doteq \sigma\{X_s \mid s \in [\tau]\}$  for  $\tau \in [n]$ . Our main object of interest for most of the chapter will be the state descriptor  $L_n(r, \tau) \doteq \sum_{i \in [n]} \mathbb{I}_{\{x_i^\tau \geq r\}}$  for any  $i \in \mathbb{N}_0, \tau \in [n]$ , which denotes the number of bins with at least  $r$  balls at time  $\tau \in [n]$ , and its expectation  $l_n(i, \tau) \doteq \mathbf{E} L_n(i, \tau)$ . The maximum bin size at time  $\tau \in [n]$  is defined as  $M_n(\tau) \doteq \min\{r \in \mathbb{N}_0 \mid L_n(r+1, \tau) = 0\}$ , and denote  $M_n^* \doteq M_n(n)$ . We define the distance  $\|x - y\| = \sum_{i=1}^n |x_i - y_i|$  for  $x = (x_1, \dots, x_n), y = (y_1, \dots, y_n) \in \mathbb{N}_0^n$ . Further a function  $f : \mathbb{N}_0^n \rightarrow \mathbb{R}$  will be called Lipschitz if  $|f(x) - f(y)| \leq \|x - y\|$ . Note that  $L_n(r, \tau) = \phi_r(X_\tau)$  where  $\phi_r(x) \doteq \sum_{i \in [n]} \mathbb{I}_{\{x_i \geq r\}}$  is a Lipschitz function.

Finally, while studying the fluid limit in Section 6.3, we will consider the time and space scaled version of the means  $g_{k,d,n}(i, t) \doteq n^{-1} l_n(i, \lfloor nt \rfloor)$  for  $i \in \mathbb{N}_0, t \in [0, 1]$ , making the dependence on  $k, d, n$  explicit. The corresponding fluid limit when  $n$  is large but  $k, d$  are fixed will be denoted by  $g_{k,d}(i, t)$ . Denote  $\tilde{d} \doteq d - k + 1$ . We will use the term asymptotically almost surely (a.a.s.) to mean that the chance that the event happens is converging to 1 as  $n \rightarrow \infty$ . The notation  $\binom{[n]}{d} \doteq \{D \mid D \subseteq [n], |D| = d\}$  will be used to subsets of  $[n]$  of size  $d$ .

## 6.2 Concentration and its consequences

### 6.2.1 Concentration

In this section we will extend the concentration results from (77) to our setup. The number of bins  $n \in \mathbb{N}$  in this sub-section will remain fixed. Recall that our system configuration at any time  $\tau \in [n]$  can be described by a vector in  $X_\tau \in \mathbb{N}_0^n$ , and that we call a function  $f : \mathbb{N}_0^n \rightarrow \mathbb{R}$  Lipschitz if for any  $x, y \in \mathbb{N}_0^n$

$$|f(x) - f(y)| \leq \|x - y\|$$

Then we may show the following concentration result.

**Lemma 8.** *Let  $f : \mathbb{N}_0^n \rightarrow \mathbb{R}$  be a Lipschitz function. Then for any time  $\tau \in [n]$  and  $y \geq 0$*

$$P(|f(X_\tau) - \mathbf{E}f(X_\tau)| \geq y) \leq 2 \exp\left(-\frac{y^2}{2k^2\tau}\right) \quad (6.1)$$

*Proof.* Since we use a tie-breaking rule, there is a unique function  $h : \mathbb{N}_0^n \times \binom{[n]}{d} \rightarrow \mathbb{N}_0^n$  so that  $X_\tau = h(X_{\tau-1}, \mathbf{D}_t)$  where  $\binom{[n]}{d} \doteq \{D \mid D \subseteq [n], |D| = d\}$  and  $\mathbf{D}_t \in \binom{[n]}{d}$  is the random choice of the  $d$  bins at time  $t$ . For any  $x, y \in \mathbb{N}_0^n$  and  $D \in \binom{[n]}{d}$ , we will now show

$$\|h(x, D) - h(y, D)\| \leq \|x - y\|. \quad (6.2)$$

Let  $I_x, I_y \subset [n]$  be the indices of the  $k$ -smallest bins among  $D$  as chosen by  $h$  for  $x$  and  $y$  respectively. Then  $|I_x| = |I_y| = k$  and

$$\|h(x, D) - h(y, D)\| - \|x - y\| = \sum_{i \in I_x \setminus I_y} \hat{\mathbb{I}}\{x_i \geq y_i\} + \sum_{j \in I_y \setminus I_x} \hat{\mathbb{I}}\{y_j \geq x_j\} \quad (6.3)$$

where  $\hat{\mathbb{I}}\{\text{condition}\}$  is 1 if condition is true, and  $-1$  if condition is false. Note that we cannot have  $x_i \geq y_i$  for some  $i \in I_x \setminus I_y$ , and  $y_j \geq x_j$  for some  $j \in I_y \setminus I_x$  simultaneously. This is because we have must  $x_j \geq x_i$ ,  $y_i \geq y_j$  because  $i \in I_x$  and  $j \in I_y$ , which would show  $x_i = x_j = y_i = y_j$ . Hence by the tie-breaking rule, one of  $i$  and  $j$  should be in  $I_x \cap I_y$ , a contradiction. Since  $|I_x \setminus I_y| = |I_y \setminus I_x| = k - |I_x \cap I_y|$ , this shows r.h.s. of (6.3) is  $\leq 0$ . This shows (6.2).

Now to finish the proof, we will use McDiarmid's finite difference inequality. For every  $\tau \in [n]$ , define  $h_\tau : \binom{[n]}{d}^\tau \rightarrow \mathbb{N}_0^n$  recursively by  $h_1(D_1) \doteq h(\mathbf{0}, D_1)$  and  $h_\tau(D_1, D_2, \dots, D_\tau) \doteq h(h_{\tau-1}(D_1, \dots, D_{\tau-1}), D_\tau)$  for  $\tau \geq 2$ . These are defined so that  $X_\tau = h_\tau(\mathbf{D}_1, \dots, \mathbf{D}_\tau)$ , where  $\mathbf{D}_1, \dots, \mathbf{D}_\tau \in \binom{[n]}{d}$  are the i.i.d. choices for the selection of bins at the  $\tau$  time points. Hence  $f(X_\tau) = f(h_\tau(\mathbf{D}_1, \dots, \mathbf{D}_\tau)) \doteq \tilde{f}(\mathbf{D}_1, \dots, \mathbf{D}_\tau)$ , and we will show that  $\tilde{f}$  satisfies the finite difference conditions.

Fix any  $\vec{D} = (D_1, \dots, D_\tau) \in \binom{[n]}{k}^\tau$  and suppose  $\vec{D}' = (D_1, \dots, D'_s, \dots, D_\tau) \in \binom{[n]}{k}^\tau$  differs only in the  $s$ th coordinate from  $\vec{D}$  for some  $1 \leq s \leq \tau$  and  $D'_s \in \binom{[n]}{d}$ . Then

$$\begin{aligned} \left| \tilde{f}(\vec{D}) - \tilde{f}(\vec{D}') \right| &= \left| f(h_t(\vec{D})) - f(h_t(\vec{D}')) \right| \leq \left\| h_t(\vec{D}) - h_t(\vec{D}') \right\| \\ &\leq \left\| h_s(D_1, \dots, D_{s-1}, D_s) - h_s(D_1, \dots, D_{s-1}, D'_s) \right\| \\ &= \left\| h(Y, D_s) - h(Y, D'_s) \right\| \leq \|h(Y, D_s) - Y\| + \|h(Y, D'_s) - Y\| = 2k, \end{aligned}$$

where we have used the Lipschitz property of  $f$  in the first line, multiple applications of the recursive definition of  $h_\tau$  and (6.2) in the second line, and the recursive definition of  $h_s$ ,  $Y \doteq h_{s-1}(D_1, \dots, D_{s-1})$ , and the property that  $h$  adds exactly  $k$  balls to the system in the third line. Hence  $\tilde{f}$  satisfies the finite difference condition. Now McDiarmid's independent bounded difference inequality (84) shows (6.1).  $\square$

For  $r \in \mathbb{N}_0$  and  $\tau \in [n]$ , recall that  $L_n(r, \tau)$  denotes the number of bins with at least  $r$  balls at time  $\tau$  and  $l_n(r, \tau) \doteq \mathbf{E}L_n(r, \tau)$ . Further, note that  $L_n(r, \tau) = \phi_r(X_\tau)$  where  $\phi_r(x) \doteq \sum_{i=1}^n \mathbb{I}_{\{x_i \geq r\}}$  is a Lipschitz function. Hence Lemma 8 has the following consequence:

**Corollary 1.** *For any  $y \geq 0$ ,  $r \in \mathbb{N}_0$ , and  $\tau \in [n]$ :*

$$\mathbf{P}(|L_n(r, \tau) - l_n(r, \tau)| \geq y) \leq 2 \exp\left(-\frac{y^2}{2k^2\tau}\right) \quad (6.4)$$

One may also show the following:

**Corollary 2.**

$$\mathbf{P}\left(\sup_{\tau \in [n]} \sup_{r \geq 0} |L_n(r, \tau) - l_n(r, \tau)| \geq k(n)\sqrt{n} \ln n\right) = \exp(-\Omega(\log^2 n)) \quad (6.5)$$

*Proof.* Note that  $L(\tau, 0) = n$  and  $L_n(\tau, n+1) = 0$  for any  $\tau \in [n]$ . Hence combining a union bound over  $r \in \{1, \dots, n\}$  along with (6.4) shows

$$\begin{aligned} \mathbf{P} \left( \sup_{\tau \in [n]} \sup_{r \geq 1} |L_n(r, \tau) - l_n(r, \tau)| \geq k(n)\sqrt{n} \ln n \right) &\leq \sum_{\tau=1}^n \sum_{r=1}^n \mathbf{P} (|L_n(r, \tau) - l_n(r, \tau)| \geq k(n)\sqrt{n} \ln n) \\ &\leq 2n^2 \exp(-\ln^2 n) = \exp(-\Omega(\ln^2 n)) \end{aligned}$$

□

Finally concentration a can also be used to approximate moments.

**Corollary 3.** *For any twice continuously differentiable function  $f : [0, 1] \rightarrow \mathbb{R}$  with  $\|f''\|_\infty \leq M$ . Then for any  $r \geq 0$  and  $\tau \in [n]$ , we have*

$$\left| \mathbf{E}f \left( \frac{L_n(r, \tau)}{n} \right) - f \left( \frac{l_n(i, \tau)}{n} \right) \right| \leq \frac{2Mk^2}{n}. \quad (6.6)$$

*Proof.* Let  $Y \doteq \frac{L_n(r, \tau)}{n}$  and  $\mu = \mathbf{E}Y = \frac{l_n(i, \tau)}{n}$ . Then by the Taylor's expansion for  $f$

$$f(Y) = f(\mu) + (Y - \mu)f'(\mu) + \int_{\mu}^Y (Y - s)f''(s)ds$$

Taking expectations and rearranging terms we obtain

$$|\mathbf{E}f(Y) - f(\mu)| = \left| \mathbf{E} \int_{\mu}^Y (Y - s)f''(s)ds \right| \leq \mathbf{E} \left| \int_{\mu}^Y (Y - s)f''(s)ds \right| \leq \frac{M}{2} \mathbf{E}(Y - \mu)^2.$$

Corollary 1 shows that  $Y$  satisfies the concentration inequality  $P(|Y - \mu| > t) \leq 2 \exp(-\frac{nt^2}{2k^2})$ . This shows  $\mathbf{E}(Y - \mu)^2 = \int_0^\infty P(|Y - \mu|^2 > t)dt \leq \frac{4k^2}{n}$ , and completes the proof using the display above. □

### 6.2.2 Maximum is concentrated around the constant $m^*(n)$

Define the constant  $m^*(n) \doteq \min \{r \mid l_n(r, n) \leq k\sqrt{n} \ln n\}$ . In this section we will show that the maximum bin size at time  $\tau = n$  will satisfy  $M_n^* \in \{m^*(n) - 1, m^*(n), m^*(n) + 1\}$  asymptotically almost surely (a.a.s) under various regimes of  $(k, d)$ . The lower bound is almost immediate, since

by Corollary 2

$$\sup_{\tau \in [n]} \sup_{r \geq 0} |L_n(r, \tau) - l_n(r, \tau)| \leq k\sqrt{n} \ln n \quad \text{a.a.s.}, \quad (6.7)$$

and hence  $L_n(m^*(n) - 1, n) > 0$  or  $M_n^* \geq m^*(n) - 1$  a.a.s. For the upper bound first we note that

$$L_n(m^*(n), n) \leq k\sqrt{n} \ln n \quad \text{a.a.s.} \quad (6.8)$$

and then we apply the following lemma to show that  $L_n(m^*(n) + 2, n) = 0$  a.a.s.

**Lemma 9.** *For any  $r, m \in \mathbb{N}$ , consider the constant*

$$p_n(m) \doteq \mathbf{P}\left(|\mathbf{D} \cap [m]| \geq \tilde{d}\right) \leq \binom{d}{\tilde{d}} \left(\frac{m+d}{n}\right)^{\tilde{d}} \mathbb{I}_{\{\tilde{d} \leq m\}} \quad (6.9)$$

where  $\tilde{d} \doteq d - k + 1$  and  $\mathbf{D} \in \binom{[n]}{d}$  is a uniformly chosen subset of  $[n]$  of size  $d$ . Then we can couple  $\{L_n(r, \tau)\}_{r \in \mathbb{N}_0, \tau \in [n]}$  and a binomial random variable  $Y = \text{Bin}(n, p_n(m))$  so that almost surely

$$L_n(r + 1, n) \mathbb{I}_{\{L_n(r, n) \leq m\}} \leq kY. \quad (6.10)$$

*Proof.* Note that  $L_n(r + 1, \tau)$  will increase at time  $\tau \in [n]$  if at least one of the  $k$  balls is assigned to a bin of size exactly  $r$ . Let  $\tilde{d} \doteq d - k + 1$ , this means that at least  $\tilde{d}$  of the  $d$  random bin choices  $\mathbf{D}_\tau \in \binom{[n]}{d}$ , at time  $\tau$ , should have size  $r$  or greater. The probability of the latter event conditioned on  $\mathcal{F}_{\tau-1}$  is  $p_n(L_n(r, \tau - 1))$ , where  $p_n$  is as defined in (6.9). Hence

$$\begin{aligned} L_n(r + 1, n) \mathbb{I}_{\{L_n(r, n) \leq m\}} &= L_n(r + 1, n) \mathbb{I}_{\{L_n(i, n) \leq m\}} \\ &\leq \sum_{\tau=1}^n (L_n(r + 1, \tau) - L_n(r + 1, \tau - 1)) \mathbb{I}_{\{L_n(r, \tau-1) \leq m\}}, \end{aligned} \quad (6.11)$$

which may be bounded by (6.10), since  $p_n(\cdot)$  is an increasing function, and each of the increments are at most  $k$ . For the bound in (6.9) see (6.18) in the next section.  $\square$

Applying Lemma 9 with  $r = m^*(n)$  and  $m = k\sqrt{n} \ln n$  along with (6.8) shows

$$L_n(m^*(n) + 1, n) \leq k \text{Bin}\left(n, \binom{d}{\tilde{d}} \left(\frac{k \ln n}{\sqrt{n}} + \frac{d}{n}\right)^{\tilde{d}} \mathbb{I}_{\{\tilde{d} \leq k\sqrt{n} \log n\}}\right) + o_P(1) \quad (6.12)$$

Now in the remainder of this section, for various regimes of  $k < d$ , we will use (6.12) to show that  $L_n(m^*(n) + 1, n) \xrightarrow{P} 0$  (or that  $L_n(m^*(n) + 2, n) \xrightarrow{P} 0$  when  $\tilde{d} = 2$ ), thereby showing that  $M_n^* \in \{m^*(n) - 1, m^*(n), m^*(n) + 1\}$  a.a.s. as  $n \rightarrow \infty$ . In the following we suppress the dependence of  $n$  in  $m^*(n)$  for brevity.

### 6.2.2.1 Regime $k < d = O(1)$

Under the assumption  $d < k$ , we have  $\tilde{d} \geq 2$ . We will consider the two cases  $\tilde{d} \geq 3$  and  $\tilde{d} = 2$  separately.

- If  $\tilde{d} \geq 3$ , then by using Markov's inequality in (6.12)

$$P(L_n(m^* + 1, n) > 0) \leq kn \binom{d}{\tilde{d}} \left( \frac{k \ln n}{\sqrt{n}} + \frac{d}{n} \right)^{\tilde{d}} + o(1) \rightarrow 0 \quad (6.13)$$

- If  $\tilde{d} = 2$ , concentration inequalities (84) for the Binomial distribution (6.12) would show that  $L_n(m^* + 1, n) \leq C \ln^2 n$  a.a.s for some constant  $C > 0$ . Applying Lemma 9 with  $r = m^* + 1$  and  $m = C \ln^2 n$  then shows that  $P(L_n(m^* + 2, n) > 0) \xrightarrow{P} 0$ .

### 6.2.2.2 Regime $k = O(n^{1/2-\epsilon}) \ll d$

Assume for some  $\epsilon \in (0, 1)$  that  $k \leq n^{1/2-\epsilon}$  and that  $k \ll d$ . In this regime we have  $d \rightarrow \infty$  and  $\tilde{d}/d \rightarrow 1$ . Now we consider the complementary cases  $d > k\sqrt{n} \ln n + k - 1$  and  $d \leq 2k\sqrt{n} \ln n$  separately.

- If  $d > k\sqrt{n} \ln n + k - 1$ , then  $\tilde{d} > k\sqrt{n} \ln n$  and, by (6.12),  $L_n(m^* + 1, n) \xrightarrow{P} 0$ .
- Assume  $d < 2k\sqrt{n} \ln n$ . Starting from (6.12), use Markov's inequality and  $\binom{d}{\tilde{d}} = \binom{d}{k-1}$ :

$$\begin{aligned} P(L_n(m^* + 1, n) > 0) &\leq nk \binom{d}{k-1} \left( \frac{3k \ln n}{\sqrt{n}} \right)^{\tilde{d}} + o(1) \leq n^{k+1} \left( \frac{3 \ln n}{n^\epsilon} \right)^{\tilde{d}} + o(1) \\ &= \exp \left( (k+1) \ln n - \tilde{d}(\epsilon \ln n - \ln(3 \ln n)) \right) + o(1). \end{aligned}$$

Note that the above probability is converging to zero since  $k \ll \tilde{d}$ .

### 6.3 Fluid limit approximation

Our aim in the remainder of this chapter will be to estimate  $m^*(n) \doteq \min\{r \in \mathbb{N}_0 \mid l_n(r, n) \leq k\sqrt{n} \log n\}$ . In this Section we will show that the time and space scaled mean function  $g_{k,d,n}(r, t) \doteq n^{-1}l_n(r, \lfloor tn \rfloor)$  for  $t \in [0, 1]$  and  $r \in \mathbb{N}_0$  can be approximated as  $n \rightarrow \infty$  by the “fluid limit”  $g_{k,d}(r, t)$ , which is the unique solution to a certain sequence of differential equations. By analyzing the fluid limit we will finally estimate  $m^*(n)$  in Section 6.4.

#### 6.3.1 Approximate recursion for the means

In Section 6.3.1, we establish that the means  $\{l_n(r, \tau)\}_{r \in \mathbb{N}_0, \tau \in [n]}$  approximately satisfy certain recursive equations.

##### 6.3.1.1 Notation and process evolution

Recall that each time  $\tau \in [n]$ , the  $d$  bins are selected uniformly at random without replacement from the  $n$  bins, and  $k$  balls are placed in the least loaded of these  $d$  selected bins. Now we will establish some notation to mathematically describe this process.

Suppose at each time  $\tau \in [n]$ , the bins are re-numbered from  $1, \dots, n$  in increasing order of their sizes so that bins that are preferred by the tie-breaking rule receive a lower number. Let  $\mathbf{D}_\tau \subset [n]$  denote the random choice of  $d$  points under this new numbering, and note that its distribution is still uniform on  $\binom{[n]}{d}$ . Finally, let  $\mathbf{K}_\tau \subset \mathbf{D}_\tau$  denote the subset of  $k$  smallest indices among  $\mathbf{D}_\tau$ . We then have the equation

$$L_n(r, \tau + 1) - L_n(r, \tau) = |\mathbf{K}_{\tau+1} \cap (n - L_n(r - 1, \tau), n - L_n(r, \tau))| \quad (6.14)$$

for any  $r \in \mathbb{N}, \tau \in [n - 1]$ , and  $L_n(0, \tau) = n$  for all  $\tau \in [n]$ . Note that the interval  $(n - L_n(r - 1, \tau), n - L_n(r, \tau)) \cap [n]$  corresponds to the numbering of balls that have size exactly  $r - 1$  at time  $\tau$ , and hence the right hand side in (6.14) exactly corresponds to the number of balls that enter into bins of size  $r - 1$  at time  $\tau + 1$ , which is same as the left hand side.

Let us simplify (6.14) further. Let  $C_n(\tau, r) \doteq (0, n - L_n(r, \tau)]$ , then

$$\begin{aligned}
L_n(r, \tau + 1) - L_n(r, \tau) &= |\mathbf{K}_{\tau+1} \cap C_n(r, \tau)| - |\mathbf{K}_{\tau+1} \cap C_n(r - 1, \tau)| \\
&= |\mathbf{D}_{\tau+1} \cap C_n(r, \tau)| \wedge k - |\mathbf{D}_{\tau+1} \cap C_n(r - 1, \tau)| \wedge k \\
&= (k - |\mathbf{D}_{\tau+1} \cap C_n(r - 1, \tau)|)^+ - (k - |\mathbf{D}_{\tau+1} \cap C_n(r, \tau)|)^+ \\
&= (|\mathbf{D}_{\tau+1} \cap (n - L_n(r - 1, \tau), n]| - l)^+ - (|\mathbf{D}_{\tau+1} \cap (n - L_n(r, \tau), n]| - l)^+
\end{aligned} \tag{6.15}$$

where in the second line we have used that  $C_n(\tau, r)$  is a downward closed set, in the third line that  $a \wedge b = b - (b - a)^+$ , and in the fourth line that  $l \doteq d - k$  and  $|\mathbf{D}_{\tau+1} \cap C| = d - |\mathbf{D}_{\tau+1} \cap \bar{C}|$ .

### 6.3.1.2 Recursion for the mean function

Now by taking expectations in (6.15) and approximating the sampling process we can show the following.

**Proposition 5.** Define the function  $f_{d,k}(p) \doteq \mathbf{E}(\text{Bin}(d, p) - d + k)^+$  for  $p \in [0, 1]$ . Then

$$l_n(r, \tau + 1) - l_n(r, \tau) = \mathbf{E}f_{d,k}\left(\frac{L_n(r - 1, \tau)}{n}\right) - \mathbf{E}f_{d,k}\left(\frac{L_n(r, \tau)}{n}\right) + \xi_n(r, \tau) \tag{6.16}$$

for some constants  $\{\xi_n(r, \tau)\}_{r \in \mathbb{N}, \tau \in [n-1]}$  that satisfy  $\sup_{r \in \mathbb{N}} \sup_{\tau \in [n-1]} |\xi_n(r, \tau)| \leq \frac{4d^2}{n}$ .

*Proof.* Recall  $\mathcal{F}_\tau$  is the sigma field of events upto time  $\tau$ . Taking conditional expectation with respect to  $\mathcal{F}_\tau$  in (6.15), we get:

$$\mathbf{E}[L_n(r, \tau + 1) - L_n(r, \tau) \mid \mathcal{F}_\tau] = h_n(L_n(r - 1, \tau)) - h_n(L_n(r, \tau)) \tag{6.17}$$

where  $h_n(m) \doteq \mathbf{E}(|\mathbf{D} \cap [m]| - d + k)^+$  and  $\mathbf{D}$  is a uniformly chosen subset of size  $d$  from  $[n]$ .

Since  $\frac{x+t}{y+t}$  is monotonically increasing in  $t$  when  $x \leq y$ , for any  $m$  we may find a coupling of  $|\mathbf{D} \cap [m]|$  and binomial random variables so that

$$\text{Bin}\left(d, \left(\frac{m-d}{n}\right)^+\right) \leq |\mathbf{D} \cap [m]| \leq \text{Bin}\left(d, \frac{m+d}{n} \wedge 1\right) \tag{6.18}$$



This shows that for all  $m \in [n]$ :

$$f_{d,k} \left( \left( \frac{m-d}{n} \right)^+ \right) \leq h_n(m) \leq f_{d,k} \left( \frac{m+d}{n} \wedge 1 \right)$$

where  $f_{d,k}$  is as defined in the statement of the proposition. Note from calculations in Section 6.5.1 that the derivative of  $f_{d,k}$ , and hence its Lipschitz coefficient, is bounded by  $d$ . Hence from (6.3.1.2) and the monotonicity of  $f_{d,k}$ , we obtain

$$\left| f_{d,k} \left( \frac{m}{n} \right) - h_n(m) \right| \leq \frac{2d^2}{n}. \quad (6.19)$$

Taking expectations in (6.17) and using (6.19) completes the proof.  $\square$

Note that (6.16) is not a recursive approximation for  $l_n$  yet, because the right side depends on the  $f_{d,k}$  moments of  $L_n$ . But by using concentration results from Section 6.2 the right side can indeed be approximated in term of  $l_n$ . Using this we can derive the following result.

**Lemma 10.** *Define the scaled function  $g_{k,d,n}(r, t) \doteq n^{-1}l_n(r, \lfloor nt \rfloor)$  for any  $r \in \mathbb{N}_0$  and  $t \in [0, 1]$ . Then there is a sequence of step functions  $\{\eta(r, \cdot)\}_{r \in \mathbb{N}}$  on  $[0, 1]$  with  $\sup_{r \in \mathbb{N}} \sup_{t \in [0, 1]} |\eta(r, t)| \leq \frac{4d^2k^2}{n}$ , such that for any  $r \in \mathbb{N}$  and  $t \in [0, 1]$*

$$g_{k,d,n}(r, t) = \int_0^{\lfloor \frac{nt}{n} \rfloor} f_{d,k}(g_{k,d,n}(r-1, s))ds - \int_0^{\lfloor \frac{nt}{n} \rfloor} f_{d,k}(g_{k,d,n}(r, s))ds + \eta(r, t). \quad (6.20)$$

Finally, note that the boundary conditions  $g_{k,d,n}(0, t) = 1$  for  $t \in [0, 1]$  and  $g_{k,d,n}(r, 0) = 0$  for every  $r \in \mathbb{N}$  are also satisfied.

*Proof.* Calculations from Section 6.5.1 show that  $f''(p) = d \text{ beta}(p; d-k, k)$  where  $\text{beta}(x; \alpha, \beta) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$  is the probability density function of the beta distribution with parameters  $\alpha, \beta$ . Hence combining Corollary 3 along with Proposition 5 we obtain

$$l_n(r, \tau+1) - l_n(r, \tau) = f_{d,k} \left( \frac{l_n(r-1, \tau)}{n} \right) - f_{d,k} \left( \frac{l_n(r, \tau)}{n} \right) + \eta'(r, \tau) \quad (6.21)$$

for some  $\eta'(\cdot, \cdot)$  that satisfy

$$\sup_{r \in \mathbb{N}} \sup_{\tau \in [n-1]} |\eta'(r, \tau)| \leq \frac{H(d, k)}{n} \quad (6.22)$$

where  $H(d, k) \doteq 4d(\sup_{p \in [0,1]} \text{beta}(p; d-k, k) k^2 + d)$ . Next for any  $a, b \in \mathbb{N}_0$  note :

$$\begin{aligned} \sup_{p \in [0,1]} \text{beta}(p; a+1, b+1) &= \text{beta}\left(\frac{a}{a+b}; a+1, b+1\right) \\ &= (a+b+1) \binom{a+b}{a} \left(\frac{a}{a+b}\right)^a \left(\frac{b}{a+b}\right)^b \\ &\leq (a+b+1) \left(\frac{a}{a+b} + \frac{b}{a+b}\right)^{a+b} = (a+b+1) \end{aligned} \quad (6.23)$$

where we use the convention that  $0! = 1$  and  $0^0 = 1$ . Taking  $a = d-k-1$  and  $b = k-1$  this shows  $H(d, k) \leq 4d^2k^2$ .

Finally, for any  $r \in \mathbb{N}$ , adding over (6.21) for  $\tau \in \{1, \dots, \lfloor nt \rfloor\}$  and dividing by  $n$  shows

$$g_{k,d,n}(r, t) = \int_0^{\lfloor nt \rfloor / n} f_{d,k}(g_{k,d,n}(r-1, s)) ds - \int_0^{\lfloor nt \rfloor / n} f_{d,k}(g_{k,d,n}(r, s)) ds + \eta(r, t) \quad (6.24)$$

where  $\eta(r, t) \doteq n^{-1} \sum_{\tau=1}^{\lfloor nt \rfloor} \eta'(r, \tau)$ . From (6.22) note that for any  $r \in \mathbb{N}$  and  $t \in [0, 1]$

$$|\eta(r, t)| \leq n^{-1} \sum_{\tau=1}^{\lfloor nt \rfloor} |\eta'(r, \tau)| \leq \frac{H(d, k) \lfloor nt \rfloor}{n^2} \leq \frac{4d^2k^2}{n}.$$

The boundary conditions immediately follow from the definition of  $g_{k,d,n}$ . □

### 6.3.2 Error estimates

Based on Lemma 10, we will now show the  $g_{k,d,n}$  is close to the fluid limit  $g_{k,d}$  defined below. Hence our aim in this subsection  $\Delta_n(i, t) \doteq |g_{k,d,n}(i, t) - g_{k,d}(i, t)|$  using Gronwall's inequality.

**Definition 7.** For any  $1 \leq k < d \in \mathbb{N}_0$  the *fluid limit* is the sequence of functions  $\{g_{k,d}(i, t)\}_{i \geq 1, t \in [0,1]}$  that is the unique solution to following equations

$$g_{k,d}(i, t) = \int_0^t f_{d,k}(g_{k,d}(i-1, s)) ds - \int_0^t f_{d,k}(g_{k,d}(i, s)) ds \quad \forall i \in \mathbb{N} \quad (6.25)$$

with the boundary  $g_{d,k}(0, t) = 1$  for  $t \in [0, 1]$ .

**Remark 6.3.1.** Since  $f_{d,k}$  is a Lipschitz function, standard theory, by induction on  $i \in nat$ , shows that the fluid limit is well defined.

Recall  $\Delta_n(i, t) \doteq |g_{k,d,n}(i, t) - g_{k,d}(i, t)|$ . Then using (6.20), (6.25), and the fact that  $f_{d,k}$  is convex and  $d$ -Lipschitz, we have

$$\Delta_n(i, t) \leq \left( d \int_0^t \Delta_n(i-1, s) ds + C \frac{d^2 k^2}{n} \right) + \int_0^t f'_{d,k}(H(i, s)) \Delta_n(i, s) ds \quad (6.26)$$

for every  $i \geq 1$ , where  $H_n(i, t) \doteq g_{k,d,n}(i, t) \vee g_{k,d}(i, t)$ . Denoting  $\Delta_n(i) \doteq \sup_{t \in [0,1]} \Delta_n(i, t)$ , Gronwall's inequality shows the recursion

$$\begin{aligned} \Delta_n(i) &\leq \left( d \Delta_n(i-1) + \frac{C d^2 k^2}{n} \right) \exp \left( \int_0^t f'_{d,k}(H_n(i, t)) ds \right) \\ &\leq \left( d \Delta_n(i-1) + \frac{C d^2 k^2}{n} \right) \exp \left( \int_0^t d \binom{d-1}{k-1} H_n(i, t)^{d-k} ds \right) \end{aligned} \quad (6.27)$$

where we have used evaluated the derivative of  $f_{d,k}$  in the last line (see Section 6.5.1).

### 6.3.2.1 Explicit calculations for the case $k = 1$

The fluid limit (6.25) for  $k = 1$  satisfies

$$g_{1,d}(i, t) = \int_0^t g_{1,d}(i-1, s)^d ds - \int_0^t g_{1,d}(i, s)^d ds \quad \forall i \in \mathbb{N}, t \in [0, 1] \quad (6.28)$$

This shows that  $g_{1,d}(i, s) \leq s$  for any  $i \geq 1$ . Similarly examining (6.17), shows that  $g_{1,d,n}(i, s) \leq s$ . Hence  $H_n(i, s) \leq s$  for  $i \geq 1$  and  $s \in [0, 1]$ . Using this in (6.27) shows for any  $i \geq 1$

$$\begin{aligned} \Delta_n(i) &\leq \left( d \Delta_n(i-1) + \frac{C d^2}{n} \right) \exp \left( \int_0^t d H_n(i, s)^{d-1} ds \right) \\ &\leq \left( d \Delta_n(i-1) + \frac{C d^2}{n} \right) \exp \left( \int_0^t d s^{d-1} ds \right) \leq A(d) \Delta_n(i-1) + B(d, n) \end{aligned} \quad (6.29)$$

where  $A(d) \doteq de$  and  $B(d, n) \doteq \frac{Ced^2}{n}$ . Expanding this recursively, we get

$$\Delta_n(i) \leq B(d, n) \frac{A(d)^i - 1}{A(d) - 1} \leq \frac{Ced^2(de)^i}{n} \quad (6.30)$$

In particular, if  $d = O(\log n)$ , then this shows that  $\|g_{1,d,n}(i, \cdot) - g_{1,d}(i, \cdot)\|_* \rightarrow 0$  for each  $i \in \mathbb{N}$  as  $n \rightarrow \infty$ . This also shows that some of the properties of  $g_{1,d,n}(i, \cdot)$  extend to  $g_{1,d}(i, \cdot)$ : namely,  $g_{1,d}(i, t) \in [0, 1]$  is non-decreasing function of  $t$  for a fixed  $i$ , and a non-increasing function of  $i$  for a

fixed  $t$ . Further, we now show that the functions  $g_{1,d}(i, \cdot)$  converge to zero uniformly as  $i \rightarrow \infty$ . Let  $h(t) \doteq \lim_{i \rightarrow \infty} g_{1,d}(i, t)$  for any  $t \in [0, 1]$  (which exists by monotonicity in  $i$ ). We may let  $i \rightarrow \infty$  in (6.28) and use dominated convergence theorem to conclude

$$h(t) = \int_0^t h(s)^d ds - \int_0^t h(s)^d ds \quad \forall t \in [0, 1]. \quad (6.31)$$

This shows  $h = 0$  and hence  $g_{1,d}(i, \cdot)$  converge uniformly to the 0 function as  $i \rightarrow \infty$ .

### 6.3.3 Properties of the fluid limit

Here we will provide upper and lower bounds for fluid limit solution from Definition 7.

#### 6.3.3.1 Bounds for the case $k = 1$

We will estimate the growth of  $g_{1,d}$  from (6.28). Define

$$u(i, d) \doteq \left( \prod_{j=1}^i \left( \sum_{k=0}^{j-1} d^k \right)^{d^{i-j}} \right)^{-1} \quad \text{for } i \geq 1 \quad (6.32)$$

We can show the following lemma.

**Lemma 11.** *Let  $g_{1,d}$  be the fluid limit from Definition 7 for the case  $k = 1$  and let  $u(i, d)$  be constants defined in (6.32) for each  $i \in \mathbb{N}$ . Then we have the bounds*

$$\frac{u(i, d)}{2^{\sum_{k=0}^{i-1} d^k}} t^{d+1} \leq g_{1,d}(i, t) \leq u(i, d) t^{\sum_{k=0}^{i-1} d^k} \quad (6.33)$$

for any  $t \in [0, 1]$  and any  $i \in \mathbb{N}$ .

*Proof.* The proof proceeds by induction on  $i \in \mathbb{N}$ . Since  $u(1, d) = 1$  and  $g_{1,d}(1, t) \leq t$ , the upper bound in (6.33) holds when  $k = 1$ . For the lower bound, note for any  $i \geq 1$  and  $t \in [0, 1]$

$$\begin{aligned} g_{1,d}(i, t) &\geq \frac{1}{2} \left( g_{1,d}(i, t) + \int_0^t g_{1,d}(i, s) ds \right) \geq \frac{1}{2} \left( g_{1,d}(i, t) + \int_0^t g_{1,d}(i, s)^d ds \right) \\ &= \frac{1}{2} \int_0^t g_{1,d}(i-1, s)^d ds \end{aligned} \quad (6.34)$$

where we have used that  $g_{1,d}(i, t) \geq \int_0^t g_{1,d}(i, s) ds$  since  $g_{1,d}(i, s)$  is non-decreasing in  $s$  and  $t \in [0, 1]$ ,  $g_{1,d}(i, s) \geq g_{1,d}(i, s)^d$  since  $g_{1,d}(i, s) \in [0, 1]$ , and finally (6.28) in the last line. Since  $g_{1,d}(0, s) = 1$ , (6.34) shows that the lower bound in (6.33) holds when  $i = 1$ .

Now we will prove the inductive step. Assume that (6.33) holds for some  $i \in \mathbb{N}$ . For the upper bounds, since  $g_{1,d}$  is always non-negative, note from (6.28) and the induction-hypothesis that for any  $t \in [0, 1]$

$$g_{1,d}(i+1, t) \leq \int_0^t g_{1,d}(i, s)^d ds \leq \int_0^t u(i, d)^d s^{d(\sum_{k=0}^{i-1} d^k)} ds = \frac{u(i, d)^d}{\sum_{k=0}^i d^k} t^{\sum_{k=0}^i d^k}.$$

Now the upper bound in (6.33) for  $i+1$  in place of  $i$  follows from the equality  $u(i+1, d) = \frac{u(i, d)^d}{\sum_{k=0}^i d^k}$ . For the lower bound in (6.33) for  $i+1$ , proceed similarly to use the lower bound from the induction-hypothesis, starting from (6.34).  $\square$

As a consequence of the above lemma, we may obtain the following bounds for  $t = 1$ .

**Corollary 4.** *For any  $i \in \mathbb{N}$  we have the bounds on  $g_{1,d}(i, t)$  at  $t = 1$ :*

$$\exp\left(-d^i \left(\frac{3 \ln d}{(d+1)^2} + \ln 2\right)\right) \leq g_{1,d}(i, 1) \leq \exp\left(-d^i \frac{\ln d}{2(1+d)^2}\right). \quad (6.35)$$

Further if  $d \geq e^4$ , one also has the following bound for any  $i \in \mathbb{N}$ :

$$\exp(-d^{i+1}) \leq g_{1,d}(i, 1) \leq \exp(-d^{i-2}). \quad (6.36)$$

*Proof.* The inequality (6.36) is a special case of (6.35) when  $d \geq e^4$ . We prove the latter, by taking  $t = 1$  in (6.33) and noting the inclusion

$$\begin{aligned} \ln u(i, d) &= -\sum_{j=2}^i d^{i-j} \ln\left(\frac{d^j - 1}{d - 1}\right) \in -\sum_{j=2}^i d^{i-j} [j - 1, j] \ln d \\ &= -d^i \ln d \sum_{j=2}^i \frac{[j - 1, j]}{d^j} \subseteq d^i \frac{\ln d}{(1 + d)^2} \left[-3, -\frac{1}{2}\right] \end{aligned}$$

where we have used  $\ln(\frac{d^j - 1}{d - 1}) \in [j - 1, j] \ln d$  in the first line, since  $d \geq 2$ .  $\square$

### 6.3.3.2 Case $k < d = O(1)$

Recall  $\tilde{d} \doteq d - k + 1 \geq 2$ . From Corollary 5 in Section 6.5.1 we know that  $p^{\tilde{d}} \leq f_{d,p}(p) \leq \binom{d}{k-1} p^{\tilde{d}}$  for any  $p \in [0, 1]$ . Let  $H \doteq \binom{d}{k-1} = O(1)$ . By (6.25) and monotonicity of  $g_{k,d}(i-1, \cdot)$

$$g_{k,d}(i, t) \leq \int_0^t f_{d,k}(g_{k,d}(i-1, s)) ds \leq \int_0^t H g_{k,d}(i-1, s)^{\tilde{d}} ds \leq H g_{k,d}(i-1, t)^{\tilde{d}} \quad (6.37)$$

for any  $t \in [0, 1]$ . We will assume that  $\lim_{i \rightarrow \infty} g_{k,d}(i, 1) = 0$  holds just like we have shown at the end of Section 6.3.2 for the case when  $k = 1$ . Hence there is an  $i_0 \in \mathbb{N}$  so that  $b_0 \doteq H g_{k,d}(i_0, 1) < 1$ . Repeated application of (6.37) gives  $\forall i > i_0$

$$g_{k,d}(i, t) \leq H^{1+\tilde{d}+\dots+\tilde{d}^{i-i_0-1}} g(i_0, t)^{\tilde{d}^{i-i_0}} \leq (H g_{k,d}(i_0, t))^{\tilde{d}^{i-i_0}} \leq b_0^{\tilde{d}^{i-i_0}} \quad (6.38)$$

This provides an upper bound for  $g_{k,d}$ . For the lower bound, rearrange (6.25) and use that  $f_{d,k}(p) \leq E \text{Bin}(d, p) = dp$  and the monotonicity of  $g_{k,d}$  to see

$$\begin{aligned} \int_0^t f_{d,k}(g_{k,d}(i-1, s)) ds &\leq g_{k,d}(i, t) + \int_0^t f_{d,k}(g_{k,d}(i, s)) ds \\ &\leq g_{k,d}(i, t) + t f_{d,k}(g_{k,d}(i, t)) \leq (1+d) g_{k,d}(i, t) \end{aligned}$$

Hence the bounds  $f_{d,k}(p) \geq p^{\tilde{d}}$  with  $c_2 \doteq \frac{1}{1+\tilde{d}}$  provides the lower bound

$$g_{k,d}(i, t) \geq c_2 \int_0^t g_{k,d}(i-1, s)^{\tilde{d}} ds. \quad (6.39)$$

Iterating (6.39) then shows  $g_{k,d}(1, t) \geq c_2 t$ ,  $g_{k,d}(2, t) \geq \frac{c_2^{1+\tilde{d}}}{\tilde{d}+1} t^{\tilde{d}+1}$ ,  $g_{k,d}(3, t) \geq \frac{c_2^{1+\tilde{d}+\tilde{d}^2}}{(\tilde{d}+1)^{\tilde{d}}(\tilde{d}^2+\tilde{d}+1)} t^{\tilde{d}^2+\tilde{d}+1}$ , and more generally

$$g_{k,d}(i, t) \geq \frac{c_2^{1+\tilde{d}+\tilde{d}^2+\dots+\tilde{d}^{i-1}}}{\prod_{j=1}^i \left( \frac{\tilde{d}^j-1}{\tilde{d}-1} \right)^{\tilde{d}^{i-j}}} t^{1+\tilde{d}+\dots+\tilde{d}^{i-1}} \quad \text{for any } i \in \mathbb{N}. \quad (6.40)$$

Hence, for any  $i \in \mathbb{N}$ , evaluating the above bound at  $t = 1$  we get

$$\begin{aligned}
g_{k,d}(i, 1) &\geq \prod_{j=1}^i \left( \frac{c_2(\tilde{d}-1)}{\tilde{d}^j - 1} \right)^{\tilde{d}^{i-j}} \geq \prod_{j=1}^i \left( \frac{c_2(\tilde{d}-1)}{\tilde{d}^j} \right)^{\tilde{d}^{i-j}} = (c_2(\tilde{d}-1))^{\frac{\tilde{d}^i - 1}{\tilde{d} - 1}} \left( \frac{1}{\tilde{d}} \right)^{\tilde{d}^i \sum_{j=1}^i j \tilde{d}^{-j}} \\
&\geq \left( c_2(\tilde{d}-1) \wedge 1 \right)^{\tilde{d}^i} \left( \frac{1}{\tilde{d}^{\sum_{j=1}^{\infty} j \tilde{d}^{-j}}} \right)^{\tilde{d}^i} \geq a_0^{\tilde{d}^i}
\end{aligned} \tag{6.41}$$

where  $a_0 \doteq \frac{c_2(\tilde{d}-1) \wedge 1}{\tilde{d}^{\frac{1}{(1-\tilde{d})^2}}} < 1$ .

**Remark 6.3.2.** Hence by combining (6.38) and (6.41) we have shown that there is a bounded sequence  $\{c_i\}_{i \in \mathbb{N}} \subset \mathbb{R}$  (with the bound dependent on  $d$ ) so that

$$g_{k,d}(i, 1) = \exp(-\tilde{d}^{i+c_i}) \quad \text{for each } i \in \mathbb{N}. \tag{6.42}$$

where recall  $\tilde{d} \doteq d - k + 1$ .

## 6.4 Estimating the expected maximum

Recall that our aim was to estimate  $m^*(n) \doteq \min\{i \in \mathbb{N}_0 \mid l_n(i, 1) \leq k\sqrt{n} \ln n\} = \min\{i \in \mathbb{N}_0 \mid g_{k,d,n}(i, 1) \leq \frac{k \ln n}{\sqrt{n}}\}$ . We will use the fluid approximation and its properties derived in the previous sections to establish this.

### 6.4.0.1 Calculations for the case $k = 1$ and $2 \leq d \leq C \log n$

**Lemma 12.** *Assume that  $k(n) = 1$  and there is a constant  $C$  so that  $2 \leq d(n) \leq C \ln n$ . Then there is an  $N_0 \in \mathbb{N}$  so that for each  $n \geq N_0$*

$$\frac{\ln \ln n - \ln 4}{\ln d(n)} - 1 \leq m^*(n) \leq \frac{\ln \ln n}{\ln d(n)} + 2.$$

*Proof.* Let  $i_u(n) \doteq \frac{\ln \ln n}{\ln d} + 2$ , and recall  $\Delta_n(i) \doteq \sup_{t \in [0,1]} |g_{1,d,n}(i, t) - g_{1,d}(i, t)|$ . Then for some constant  $C_1$  we have by using (6.30) and (6.36):

$$\begin{aligned} g_{1,d,n}(i_u(n), 1) &\leq g_{1,d}(i_u(n), 1) + \Delta_n(i_u(n)) \leq \exp(-d^{\frac{\ln \ln n}{\ln d}}) + \frac{C_1^{i_u(n)} (\ln n)^{2+i_u(n)}}{n} \\ &= \frac{1}{n} + \exp(i_u(n) \ln C_1 + (2 + i_u(n)) \ln \ln n - \ln n) \\ &\leq \frac{1}{n} + \exp(-\frac{1}{2} \ln n) \leq \frac{\ln n}{\sqrt{n}} \end{aligned} \quad (6.43)$$

when  $n$  is large enough. This shows that  $m^*(n) \leq i_u(n) = \frac{\ln \ln n}{\ln d} + 2$  eventually. Similarly, we will prove a lower bound. Let  $i_l(n) = \frac{\ln \ln n - \ln 4}{\ln d} - 1$  then, for any  $i \leq i_l(n)$ , using the lower bounds for (6.30) and (6.36) there are constant  $C_1$  and  $C_2$  so that

$$\begin{aligned} g_{1,d,n}(i, 1) &\geq g_{1,d}(i, 1) - \Delta_n(i) \geq \exp(-d^{i+1}) - C_2 \exp(i \ln C_1 + (2 + i) \ln \ln n - \ln n) \\ &\geq \exp(-d^{i_l(n)+1}) - C_2 \exp(i_l(n) \ln C_1 + (2 + i_l(n)) \ln \ln n - \ln n) \\ &\geq \frac{1}{\sqrt[4]{n}} - \frac{C_2}{\sqrt{n}} > \frac{\ln n}{\sqrt{n}} \end{aligned} \quad (6.44)$$

when  $n$  is large enough. This shows that  $m^*(n) \geq i_l(n) = \frac{\ln \ln n - \ln 4}{\ln d} - 1$  eventually.  $\square$

#### 6.4.0.2 Calculations for the case $1 \leq k < d = O(1)$

**Lemma 13.** *Assume that  $1 \leq k < d \leq C$  for each  $n$  for some constant  $C$ . Then there are constants  $\tilde{C}, N_0 > 0$  so that*

$$\left| m^*(n) - \frac{\ln \ln n}{\ln(d - k + 1)} \right| \leq \tilde{C} \quad (6.45)$$

for each  $n \geq N_0$ .

*Proof.* Since  $d$  is bounded, using Remark 6.3.2, we may find a constant  $K > 0$  (independent of  $n$ ) so that  $\sup_{i \in \mathbb{N}} |c_i| \leq K$  and  $g_{k,d}(i, 1) = \exp(-\tilde{d}^{i+c_i})$  for each  $i \in \mathbb{N}$ , where  $\tilde{d} \doteq d - k + 1$ . Further  $\Delta_n(i) \doteq \sup_{t \in [0,1]} |g_{k,d,n}(i, t) - g_{k,d}(i, t)|$  from (6.27) satisfies

$$\Delta_n(i) \leq A \Delta_n(i - 1) + \frac{B}{n} \quad \forall i \in \mathbb{N}. \quad (6.46)$$



for some constants  $A, B > 1$  that are independent of  $n$ . Since  $\Delta_n(0) = 0$ , this shows for some constant  $c > 0$  that  $\Delta_n(i) \leq \frac{cA^i}{n}$  for each  $i \geq 1$ . Hence with  $i_u(n) \doteq \frac{\log \log n}{\log d} + K$ ,

$$\begin{aligned} g_{k,d,n}(i_u, 1) &\leq g_{k,d}(i_u, 1) + \Delta_n(i_u) \leq \exp(-\tilde{d}^{\frac{\log \log n}{\log d}}) + \frac{cA^{i_u(n)}}{n} \\ &= \frac{1}{n} + \frac{cA^K(\log n)^{\log_{\tilde{d}} A}}{n}. \end{aligned}$$

Therefore  $g_{k,d,n}(i_1, 1) \leq \frac{k \log n}{\sqrt{n}}$  when  $n$  is large. This shows  $m^*(n) \leq i_u(n) = \frac{\log \log n}{\log d} + K$  for large  $n$ . For the lower bound let  $i_l(n) \doteq \frac{\log \log n}{\log d} - K - 2$ . Then for any  $i \leq i_l$

$$\begin{aligned} g_{k,d,n}(i, 1) &\geq g_{k,d}(i, 1) - \Delta_n(i) \geq \exp(-\tilde{d}^{i+c_i}) - \frac{cA^i}{n} \geq \exp(-\tilde{d}^{i_l(n)+c_i}) - \frac{cA^{i_l(n)}}{n} \\ &\geq \exp\left(-\tilde{d}^{\frac{\log \log n}{\log d}-2}\right) - \frac{cA^{\frac{\log \log n}{\log d}}}{n} = \frac{1}{n^{1/\tilde{d}^2}} - \frac{c(\log n)^{\log_{\tilde{d}} A}}{n}. \end{aligned}$$

Since  $\tilde{d} \geq 2$ , for large  $n$  we must have  $g_{k,d,n}(i, 1) > \frac{k \ln n}{\sqrt{n}}$  for each  $i \leq i_l$ . This shows  $m^*(n) > i_l(n) = \frac{\log \log n}{\log d} - K - 2$ .  $\square$

## 6.5 Technical estimates

In this section we present some technical estimates that were used in the previous analysis.

### 6.5.1 Properties of $f_{d,k}$

Let  $U_1, U_2 \dots U_d$  be i.i.d.  $U[0, 1]$  random variables. Then  $B(p) \doteq \sum_{i=1}^d \mathbb{I}_{\{U_i \leq p\}}$  has the binomial distribution  $\text{Bin}(d, p)$ . In this section, we will study the function  $f_{d,k}(p) \doteq \mathbf{E}(\text{Bin}(d, p) - d + k)^+ = \mathbf{E}(B(p) - d + k)^+$  for  $p \in [0, 1]$ . This is the same function that is studied in (70), but we will derive some of its properties using the probabilistic definition. Recall the constant  $\tilde{d} \doteq d - k + 1$  that will be important in many of our calculations.

**Proposition 6.**  $f_{d,k}(p)$  is monotonically increasing in  $p$ ,  $f_{d,k}(0) = 0$ , and  $f_{d,k}(1) = k$ . In fact,  $f_{d,k}(p)$  is a polynomial of degree  $\tilde{d}$  and satisfies

$$hp^{\tilde{d}} \leq f_{d,k}(p) \leq Hp^{\tilde{d}} \tag{6.47}$$

where  $h \doteq \min_{i \in \{0, \dots, k-1\}} \frac{(i+1)\binom{d}{d+i}}{\binom{k-1}{i}}$  and  $H \doteq \max_{i \in \{0, \dots, k-1\}} \frac{(i+1)\binom{d}{d+i}}{\binom{k-1}{i}}$ . Further  $h \geq 1$ .

*Proof.* Expansion of the binomial distribution shows

$$f_{d,k}(p) = \sum_{j=d-k+1}^d (j-d+k) \binom{d}{j} p^j (1-p)^{d-j} = p^{\tilde{d}} \left( \sum_{i=0}^{k-1} (i+1) \binom{d}{\tilde{d}+i} p^i (1-p)^{k-1-i} \right)$$

From the definition of  $h$  and  $H$  and binomial expansion one may see that (6.47) is satisfied. Finally note that  $\frac{(i+1)\binom{d}{d+i}}{\binom{k-1}{i}} \geq 1$  for any  $i \in \{0, \dots, k-1\}$  and hence  $h \geq 1$ .  $\square$

**Lemma 14.** *The first and second derivatives of  $f_{d,k}$  may be described in terms of the beta distribution. For any  $p \in [0, 1]$  we have*

$$f'_{d,k}(p) = d\mathbf{P}(\mathcal{B}e(d-k, k) \leq p) \quad (6.48)$$

and

$$f''_{d,k}(p) = d \text{beta}(p; d-k, k) \quad (6.49)$$

where  $\mathcal{B}e(\alpha, \beta)$  denotes a beta random variable with density  $\text{beta}(x; \alpha, \beta) \doteq \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$  for  $x \in [0, 1]$ .

*Proof.* Let  $\mathcal{I}(p) \doteq \{i \mid 1 \leq i \leq d, U_i \leq p\}$  and  $\mathcal{G}_p \doteq \sigma\{\mathcal{I}(q) \mid q \geq p\}$ , so that  $B(p) = |\mathcal{I}(p)|$  is  $\mathcal{G}_p$  measurable.

Next note for any  $p \in (0, 1)$  and small  $h > 0$

$$\begin{aligned} \frac{f_{d,k}(p) - f_{d,k}(p-h)}{h} &= \mathbf{E} \frac{(B(p) - d + k)^+ - (B(p-h) - d + k)^+}{h} \\ &= \mathbf{E} \frac{\mathbb{I}_{\{B(p) \geq \tilde{d}\}} (B(p) - B(p-h)) \wedge (B(p) - \tilde{d} + 1)}{h} \\ &= \mathbf{E} \frac{\left[ |\mathcal{I}(p) \setminus \mathcal{I}(p-h)| \wedge (|\mathcal{I}(p)| - \tilde{d} + 1) \mid \mathcal{G}_p \right]}{h} \mathbb{I}_{\{|\mathcal{I}(p)| \geq \tilde{d}\}} \end{aligned}$$

By independence of  $U_i$ ,  $|\mathcal{I}(p) \setminus \mathcal{I}(p-h)|$  is distributed as  $\text{Bin}(|\mathcal{I}(p)|, \frac{h}{p})$  when conditioned on  $\mathcal{G}_p$ .

Hence we have

$$= \mathbf{E} \frac{\left[ \text{Bin}(B(p), \frac{h}{p}) \wedge (B(p) - \tilde{d} + 1) \mid B(p) \right]}{h} \mathbb{I}_{\{B(p) \geq \tilde{d}\}}$$

Let  $h \rightarrow 0$  and use the fact that  $\frac{\mathbf{E} \text{Bin}(a, \frac{h}{p})^{\wedge r}}{h} \rightarrow \frac{a}{p}$  for any  $r \geq 1$ . Then by dominated convergence:

$$(f_{d,k}(p))' = \frac{\mathbf{E} B(p) \mathbb{I}_{\{B(p) \geq \tilde{d}\}}}{p} \quad (6.50)$$

If  $g(p) = \mathbf{E} B(p) \mathbb{I}_{\{B(p) \geq p\}}$ , then by a similar argument we can show  $g'(p) = \frac{1}{p} \left[ (\tilde{d} - 1) \tilde{d} \mathbf{P}(B(p) = \tilde{d}) + g(p) \right]$ . Hence, using the quotient rule in (6.50), this shows

$$\begin{aligned} (f_{d,k}(p))'' &= \frac{\tilde{d}(\tilde{d} - 1) \mathbf{P}(B(p) = \tilde{d})}{p^2} \\ &= \tilde{d}(\tilde{d} - 1) \binom{d}{\tilde{d}} p^{\tilde{d}-2} (1-p)^{d-\tilde{d}} \\ &= d \text{beta}(p; \tilde{d} - 1, d - \tilde{d} + 1) = d \text{beta}(p; d - k, k) \end{aligned} \quad (6.51)$$

This shows (6.49). Now integrating (6.49) shows (6.48).  $\square$

**Remark 6.5.1.** The above Lemma shows that  $f'_{d,k}(p)$  is a increasing function of  $p$  that is always bounded by  $d$ . This shows that  $f_{d,k}(p)$  is convex function with Lipschitz constant  $d$ .

**Corollary 5.** *Integrating (6.51) derived in the proof of the above lemma shows that  $f_{d,k}(p)' \leq \tilde{d} \binom{d}{\tilde{d}} p^{\tilde{d}-1}$  and  $f_{d,k}(p) \leq \binom{d}{\tilde{d}} p^{\tilde{d}}$ . Hence combining with Proposition 6 we see that*

$$p^{\tilde{d}} \leq f_{d,k}(p) \leq p^{\tilde{d}} \binom{d}{k-1} \quad \text{for any } p \in [0, 1].$$

## CHAPTER 7

### Limit theorems for the Supermarket model

#### 7.1 Introduction

Recall from Section 5.2, that the Supermarket model with  $n$ -servers and arrival rate  $\lambda_n$  under the  $JSQ(d_n)$  scheme is a processing system with  $n$  parallel queues in which each queue's jobs are processed by the associated server at rate 1. Jobs arrive at rate  $n\lambda_n$  and join the shortest queue amongst  $d_n$  randomly selected queues (without replacement), with  $d_n \in [n]$ . The interarrival times and service times are mutually independent exponential random variables. In this chapter we will discuss limit theorems for the Supermarket model under the assumption that  $d_n \rightarrow \infty$ . The  $l_1^\downarrow$  valued infinite dimensional vector  $\mathbf{G}_n(t) = (G_{n,i}(t))_{i \geq 1}$ , where  $G_{n,i}(t)$  is the fraction of queues with at least  $i$  jobs at time  $t$ , plays an important role in the analysis of Supermarket model, since  $\mathbf{G}_n$  forms a  $l_1^\downarrow$  valued continuous time Markov chain (CTMC). First, we will show a law of large numbers result that says that  $\mathbf{G}_n$  converges in probability, as  $n \rightarrow \infty$ , to an universal  $l_1^\downarrow$  valued deterministic trajectory as long as  $d_n \rightarrow \infty$ . We indirectly showed a similar result for the balls and bins problem in Chapter 6 by using concentration results and analyzing the mean trajectory of the system.

However our main aim in the chapter will be to develop diffusion approximations for  $\mathbf{G}_n$  in the critical regime (i.e. when  $\lambda_n \rightarrow 1$  in a suitable manner) that allow for possibly a slower growth of  $d_n$  than that permitted by the results in Mukherjee et al. (91). In fact the results we establish will allow for  $d_n \rightarrow \infty$  in an arbitrary manner and will recover the results of (91) in the special case  $\frac{d_n}{\sqrt{n \log n}} \rightarrow \infty$  (with a different proof). In order to motivate the type of limit theorems we seek, we begin by observing that the centering  $\mathbf{e}_1$  used in the definition of  $\mathbf{Y}_n = \sqrt{n}(\mathbf{G}_n - \mathbf{e}_1)$  used in (91, 45) is a stationary point of the fluid limit given in (7.3) with  $\lambda = 1$  and thus the results of (45) and (91) give information on fluctuations of the state process  $\mathbf{G}_n$  about this stationary point. However  $\mathbf{e}_1$  is not the only stationary point of (7.3) (when  $\lambda = 1$ ) and in fact this ODE has uncountably

many fixed points with a typical such point given as  $\mathbf{f}_k^\gamma \doteq \sum_{j=1}^k \mathbf{e}_j + \gamma \mathbf{e}_{k+1}$ , where  $\mathbf{e}_j$  is the  $j$ -th unit vector in  $l_2$  (with 1 at the  $j$ -th coordinate and zeroes elsewhere),  $k \in \mathbb{N}$  and  $\gamma \in [0, 1)$ . All of these stationary points arise in a natural fashion. Indeed, it turns out that the evolution of the state process  $\mathbf{G}_n$  can be described via the equation (see Remark 7.3.1)

$$\mathbf{G}_n(t) = \mathbf{G}_n(0) + \int_0^t [\mathbf{a}_n(\mathbf{G}_n(s)) - \mathbf{b}(\mathbf{G}_n(s))] ds + \mathbf{M}_n(t),$$

where  $\mathbf{M}_n$  is a (infinite dimensional) martingale converging to 0 in probability (see Lemma 15, Chapter 7) and  $\mathbf{a}_n, \mathbf{b}$  are certain maps from  $l_1^\downarrow$  to  $l_1$  (see Remark 7.3.1 for details). Thus for large  $n$ , trajectories of  $\mathbf{G}_n$  will be close to solutions of the infinite dimensional ODE

$$\dot{\mathbf{g}}_n = \mathbf{a}_n(\mathbf{g}_n) - \mathbf{b}(\mathbf{g}_n).$$

This equation has a unique stationary point  $\boldsymbol{\mu}_n$  which is introduced in Definition 9. The fixed point  $\boldsymbol{\mu}_n$  corresponds to the point in the state space  $l_1^\downarrow$  at which the inflow rate equals the outflow rate in the  $n$ -th system and thus it is of interest to explore system behavior in the neighborhood of this point. Since  $\mathbf{G}_n$  is approximated by  $\mathbf{g}_n$  (over any compact time interval), one can loosely interpret  $\boldsymbol{\mu}_n$  as a *near fixed point* of the state process  $\mathbf{G}_n$ . Furthermore, it can be shown (see Remark 7.2.3(iv)) that, if  $d_n \rightarrow \infty$  and  $\lambda_n \rightarrow 1$  in a suitable manner,  $\boldsymbol{\mu}_n$  can converge to any specified fixed point  $\mathbf{f}_k^\gamma$  of (7.3) and thus every fixed point of (7.3) arises from  $\boldsymbol{\mu}_n$  in a suitable asymptotic regime. In order to explore fluctuations of  $\mathbf{G}_n$  close to different fixed points of (7.3) it is then natural to study the asymptotic behavior of

$$\mathbf{Z}_n(t) \doteq \sqrt{n}(\mathbf{G}_n(t) - \boldsymbol{\mu}_n), \quad t \geq 0. \quad (7.1)$$

We note that in the regime considered in (91) where  $\frac{d_n}{\sqrt{n} \log n} \rightarrow \infty$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \alpha > 0$ ,  $\sqrt{n}(\mathbf{e}_1 - \boldsymbol{\mu}_n) \rightarrow \alpha$  and so in this case the asymptotic behavior of  $\mathbf{Z}_n$  can be read off from that of  $\mathbf{Y}_n$  (see Corollary 6 and Remark 7.2.5(v)). However in general  $\sqrt{n}(\mathbf{e}_1 - \boldsymbol{\mu}_n)$  (and more generally,  $\sqrt{n}(\mathbf{f}_k^\gamma - \boldsymbol{\mu}_n)$ ) may not be bounded and so the asymptotic behavior of  $\mathbf{Z}_n$  and  $\mathbf{Y}_n$  may be very different.

Here we obtain limit theorems for  $\mathbf{Z}_n$  as  $d_n \rightarrow \infty$  in an arbitrary fashion and  $\lambda_n \rightarrow 1$  in a suitable manner. Specifically in Theorems 9, 10 and 11 we consider the three cases:

(a)  $d_n/\sqrt{n} \rightarrow 0$ ; (b)  $d_n/\sqrt{n} \rightarrow c \in (0, \infty)$  and, (c)  $d_n/\sqrt{n} \rightarrow \infty$ , respectively. In all three regimes we consider initial conditions  $\mathbf{G}_n(0)$  such that for some  $r \in \mathbb{N}$ ,  $G_{n,m}(0) = \mu_{n,m} + o_p(n^{-1/2})$  for all  $m > r$  and in each case (under conditions on  $\lambda_n$ ) we obtain a limit process driven by a one dimensional Brownian motion with continuous sample paths in  $l_2$  which has all but finitely many coordinates 0. In particular, when  $r = 2$  in the second and the third case and  $r = k + 2$  for some  $k \in \mathbb{N}$  in the first case (and  $d_n, \lambda_n$  depend on  $k$  in a suitable fashion), one can describe the limit through a two dimensional diffusion driven by a one dimensional Brownian motion. The form of this two dimensional process in the three regimes is quite different; in the first case we get a linear diffusion (i.e. the drift is of the form  $b(y) = Ay$  for,  $y \in \mathbb{R}^2$  and some  $2 \times 2$  matrix  $A$ ); in the second case we get a diffusion with an exponential drift; and in the third case we obtain a reflected diffusion in the half space  $(-\infty, \alpha] \times \mathbb{R}$  for some  $\alpha \geq 0$ .

Although the limit processes in Theorems 9 and 10 are quite different from those obtained in (44) and (91), the limit in Theorem 11 has a similar form (in that it is a reflected diffusion in a half space) as in the above papers. However here as well there are some differences. In particular, depending on how  $\lambda_n$  approaches 1, the reflection occurs at a different barrier  $\alpha \in (0, \infty)$ ; in fact  $\alpha = \infty$  is possible as well in which case there is no reflection. Furthermore, recall that  $\mathbf{Z}_n$  is defined by centering about  $\boldsymbol{\mu}_n$ . In general  $\sqrt{n}(\boldsymbol{\mu}_n - \mathbf{e}_1)$  will diverge and thus the process  $\mathbf{Y}_n$  considered in the above cited papers may not converge in this regime. However, as noted previously, when  $d_n$  grows sufficiently fast, namely  $\frac{d_n}{\sqrt{n \log n}} \rightarrow \infty$  the process  $\mathbf{Y}_n$  will indeed converge and in that case we recover the result in (91) (in fact a slight strengthening in that the drift parameter in Corollary 6 is allowed to be 0). In addition Theorem 11 also covers the case  $\frac{d_n}{\sqrt{n \log n}} \rightarrow c \in (0, \infty)$  and situations where  $\lambda_n = 1 + O(n^{-1/2})$  (see Remark 7.2.5 (iv)). In such settings, once more both  $\mathbf{Z}_n$  and  $\mathbf{Y}_n$  converge and the limit of the latter has the same form as in (45, 91).

As is observed in Remarks 7.2.4 and 7.2.5, under conditions of Theorem 10 or Theorem 11,  $\boldsymbol{\mu}_n$  must converge to the fixed point  $\mathbf{e}_1 = \mathbf{f}_1^0$ . In contrast, Theorem 9 allows for a range of asymptotic behavior for  $\boldsymbol{\mu}_n$ . In particular, under the conditions of the theorem, with suitable  $\lambda_n, d_n$ ,  $\boldsymbol{\mu}_n$  can converge to the fixed point  $\mathbf{f}_k^0$  for an arbitrary  $k \in \mathbb{N}$  (see (24) for a similar observation). In such a setting the first  $k - 1$  coordinates of the limit process are essentially 0 (see Theorem 9 for a precise

statement) and the  $k$ -th coordinate is the first one to exhibit stochastic variability. Thus a rather novel asymptotic behavior for the  $JSQ(d_n)$  system emerges when  $d_n$  approaches  $\infty$  at significantly slower rates than those considered in (91) and  $\lambda_n$  approach 1 in a suitable manner (in relation to  $d_n$ ).

### 7.1.1 Organization and technique overview

Section 7.2 contains all our main results. The remaining Sections starting with Section 7.3 contain proofs of the main results.

We now make some comments on the proofs of Theorems 9 - 11. The starting point is a convenient semimartingale representation for the centered state process  $\mathbf{Z}_n$  in (7.43). In the study of the behavior of the drift term in this decomposition, an important ingredient is an analysis of the asymptotic properties of the near fixed point  $\boldsymbol{\mu}_n$ , and the asymptotic behavior of the function  $\beta_n$  (see Definition 8) in  $O(n^{-\frac{1}{2}})$  sized neighborhoods around the coordinates of  $\boldsymbol{\mu}_n$ . This behavior, which is different in the three regimes considered above, determines the asymptotics of the drift  $\mathbf{A}_n(\mathbf{Z}_n(s)) - \mathbf{b}(\mathbf{Z}_n(s))$ . Properties of  $\boldsymbol{\mu}_n$  are also key in arguing that, in all three cases, under our conditions,  $(Z_{n,r+1}, \dots)$  converges to 0 in probability in  $\mathbb{D}([0, \infty) : l_2)$  (see Lemma 29). The rest of the work is in characterizing the asymptotics of the finite dimensional process  $(Z_{n,1}, \dots, Z_{n,r})$ . For this study, the three regimes require different approaches. In particular, Theorem 9 hinges on a detailed understanding of the asymptotic behavior of a tridiagonal matrix function (see e.g. Lemmas 36 and 37); Theorem 10 requires an analysis of a stochastic differential equation with an exponential drift term (in particular the drift does not satisfy the usual growth conditions); and Theorem 11 is based on a careful study of excursions of the prelimit processes above the limiting reflecting barrier and properties of Skorohod maps in order to characterize the reflection properties of the limit process.

### 7.1.2 Notation and setup

For  $m \geq 1$ , let  $[m] \doteq \{1, 2, \dots, m\}$ . We will denote finite-dimensional vectors in  $\mathbb{R}^m$  as  $\vec{x}, \vec{y}$ , etc. and  $\langle \vec{x}, \vec{y} \rangle$  will denote the standard inner-product. The standard basis vectors in  $\mathbb{R}^m$  will be denoted by  $\vec{e}_i$  for  $i = 1, 2 \dots m$ . Also,  $\|\vec{x}\| \doteq \langle \vec{x}, \vec{x} \rangle$  will denote the usual Euclidean norm.

We will often use bold symbols such as  $\mathbf{x} := (x_1, x_2, \dots)$  to denote a infinite dimensional vector or function. For  $p \in \{1, 2, \dots, \infty\}$ , let  $\|\cdot\|_p$  denote the usual  $p$ -th norm on the space of infinite sequences and  $l_p \doteq \{\mathbf{x} \in \mathbb{R}^\infty \mid \|\mathbf{x}\|_p < \infty\}$ . Let  $l_1^\perp$  be as in (5.1), which is a Polish space under  $\|\cdot\|_1$ . For  $k \in \mathbb{N}$ , let  $\mathbf{f}_k = (1, 1, \dots, 1, 0, 0 \dots) \in l_1^\perp$  denote the vector with first  $k$  indices equal to 1, and  $\mathbf{e}_k = (0, \dots, 0, 1, 0 \dots) \in l_1$  denote the vector with 1 in the  $k$ th coordinate. Finally, for  $\mathbf{z} = (z_1, z_2, \dots) \in \mathbb{R}^\infty$  and  $r \in \mathbb{N}$ , let  $\mathbf{z}_{r+} \doteq (z_{r+1}, z_{r+2}, \dots) \in \mathbb{R}^\infty$  denote the vector shifted by  $r$  steps. Similar notation will be used for functions and processes with values in  $\mathbb{R}^\infty$ .

For a Polish space  $\mathbb{S}$  and  $T > 0$ , denote by  $\mathbb{C}([0, T] : \mathbb{S})$  (resp.  $\mathbb{D}([0, T] : \mathbb{S})$ ) the space of continuous functions (resp. right continuous functions with left limits) from  $[0, T]$  to  $\mathbb{S}$ , endowed with the uniform topology (resp. Skorokhod topology). Spaces  $\mathbb{C}([0, \infty) : \mathbb{S}), \mathbb{D}([0, \infty) : \mathbb{S})$  are defined similarly. For  $f \in \mathbb{D}([0, T] : \mathbb{R})$  and  $t \leq T$ , let  $|f|_{*,t} \doteq \sup_{s \in [0,t]} |f(s)|$ . Similarly for  $\mathbf{g} \in \mathbb{D}([0, T] : l_p)$ , let  $\|\mathbf{g}\|_{p,t} \doteq \sup_{s \in [0,t]} \|\mathbf{g}(s)\|_p$ .

We will use  $\mathbb{I}_{\{cond\}}$  to denote the indicator function that takes the value 1 if *cond* is true, otherwise it takes the value 0. We will denote by  $\text{id}$  the identity map,  $\text{id}(t) = t$ , on  $[0, T]$  or  $[0, \infty)$ .

We use  $\mathbf{P}$  and  $\mathbf{E}$  to denote the probability and expectation operators, respectively. For  $x, y \in \mathbb{R}$ ,  $x \wedge y$  denotes the minimum and  $x \vee y$  the maximum of  $x$  and  $y$  respectively. For any  $x \in \mathbb{R}$ ,  $x^+ = x \vee 0$  and  $x^- = (-x) \vee 0$ . We use  $\xrightarrow{P}$  and  $\Rightarrow$  to denote convergence in probability and convergence in distribution respectively on an appropriate Polish space which will depend on the context. For a sequence of real valued random variables  $(X_n, n \geq 1)$ , we write  $X_n = o_P(b_n)$  when  $|X_n|/b_n \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . For non-negative functions  $f(\cdot), g(\cdot)$ , we write  $f(n) = O(g(n))$  when  $f(n)/g(n)$  is uniformly bounded, and  $f(n) = o(g(n))$  (or  $f(n) \ll g(n)$ ) when  $\lim_{n \rightarrow \infty} f(n)/g(n) = 0$ . We write  $f(n) \sim g(n)$  if  $f(n)/g(n) \rightarrow 1$  as  $n \rightarrow \infty$ . We will use the notation  $\lambda_n \nearrow 1$  to mean that  $\lambda_n < 1$  for every  $n$  and  $\lambda_n \rightarrow 1$  as  $n \rightarrow \infty$ .

## 7.2 Main Results

Recall the process  $\mathbf{G}_n$  from Section 7.1. Our first result gives a law of large numbers (LLN) for the process  $\mathbf{G}_n$  as  $n \rightarrow \infty$ . In order to state this result we begin by recalling the one dimensional Skorohod map (cf. (63, Section 3.6.C)) with a reflecting barrier at  $\alpha \in \mathbb{R}$ . For  $\alpha \in \mathbb{R}$  and



$f \in \mathbb{D}([0, \infty) : \mathbb{R})$  with  $f(0) \leq \alpha$ , define  $\Gamma_\alpha(f), \hat{\Gamma}_\alpha(f) \in \mathbb{D}([0, \infty) : \mathbb{R})$  as

$$\Gamma_\alpha(f)(t) = f(t) - \sup_{s \in [0, t]} (f(s) - \alpha)^+, \quad \hat{\Gamma}_\alpha(f)(t) = \sup_{s \in [0, t]} (f(s) - \alpha)^+. \quad (7.2)$$

The map  $\Gamma_\alpha$  (and sometimes the pair  $(\Gamma_\alpha, \hat{\Gamma}_\alpha)$ ) is referred to as the one-dimensional Skorohod map (with reflection at  $\alpha$ ). The following wellposedness result, which is proved in Section 7.4, will be used to characterize the LLN limit of  $\mathbf{G}_n$ .

**Proposition 7.** Fix  $\mathbf{r} \in l_1^\downarrow$ . Then there is a unique  $(\mathbf{g}, \mathbf{v}) \in C([0, \infty) : l_1^\downarrow \times l_\infty)$  that solves the following system of equations

$$\begin{aligned} g_i(t) &= \Gamma_1 \left( r_i - \int_0^\cdot (g_i(s) - g_{i+1}(s)) ds + v_{i-1}(\cdot) \right)(t) \quad \forall i \geq 1, t \geq 0 \\ v_i(t) &= \hat{\Gamma}_1 \left( r_i - \int_0^\cdot (g_i(s) - g_{i+1}(s)) ds + v_{i-1}(\cdot) \right)(t) \quad \forall i \geq 1, \quad v_0(t) = \lambda t, \quad t \geq 0. \end{aligned} \quad (7.3)$$

**Remark 7.2.1.** Using the well known characterization of a one-dimensional Skorohod map, one can alternatively characterize  $(\mathbf{g}, \mathbf{v})$  as the unique pair in  $C([0, \infty) : l_1^\downarrow \times l_\infty)$  such that  $v_i$  is non-decreasing,

$$\left. \begin{aligned} g_i(t) &= r_i - \int_0^t (g_i(s) - g_{i+1}(s)) ds + v_{i-1}(t) - v_i(t) \\ v_i(t) &\geq 0, g_i(t) \leq 1, \int_0^t (1 - g_i(s)) dv_i(s) = 0 \end{aligned} \right\} \forall i \geq 1 \quad (7.4)$$

and  $v_0(t) = \lambda t$ , for all  $t > 0$  and  $v_i(0) = 0$  for all  $i \geq 0$ .

We can now present the LLN result. The proof is given in Section 7.4.

**Theorem 8.** Let  $\mathbf{r} \in l_1^\downarrow$ . Suppose that  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{r}$  in  $l_1^\downarrow$ ,  $\lambda_n \rightarrow \lambda$  and  $d_n \rightarrow \infty$ , as  $n \rightarrow \infty$ . Then  $\mathbf{G}_n \rightarrow \mathbf{g}$  in probability in  $\mathbb{D}([0, \infty) : l_1^\downarrow)$  as  $n \rightarrow \infty$ , where  $(\mathbf{g}, \mathbf{v}) \in C([0, \infty) : l_1^\downarrow \times l_\infty)$  is the unique solution of (7.3).

**Remark 7.2.2.** Note that Theorem 8 allows  $d_n \rightarrow \infty$  in an arbitrary manner. In (91, Theorem 1) it is shown that, under the assumptions of Theorem 8,  $\mathbf{G}_n$  is a tight sequence of  $\mathbb{D}([0, \infty) : l_1^\downarrow)$  valued random variables and that every subsequential weak limit  $\hat{\mathbf{g}}$  satisfies a system of equations

given as

$$\hat{g}_i(t) = r_i - \int_0^t (\hat{g}_i(s) - \hat{g}_{i+1}(s))ds + \int_0^t p_{i-1}(\hat{\mathbf{g}}(s))ds \quad \text{for } i \geq 1 \quad (7.5)$$

where

$$p_j(\hat{\mathbf{g}}(s)) = \begin{cases} \lambda - (\lambda - 1 + \hat{g}_{j+2}(s))^+ & \text{if } j = m(\hat{\mathbf{g}}(s)) - 1 \\ (\lambda - 1 + \hat{g}_{j+1}(s))^+ & \text{if } j = m(\hat{\mathbf{g}}(s)) > 0 \\ \lambda & \text{if } j = m(\hat{\mathbf{g}}(s)) = 0 \\ 0 & \text{otherwise,} \end{cases} \quad (7.6)$$

and for  $\mathbf{x} \in l_1^\downarrow$ ,  $m(\mathbf{x}) \doteq \inf\{i \mid x_{i+1} < 1\}$ . (Note that  $m(\mathbf{G}_n(t))$  is the length of the smallest queue at time  $t$ .) The uniqueness of the above system of equations was not shown in (91).

From (7.3) and the definition in (7.2) it follows that each  $v_i$  is absolutely continuous and, for a.e.  $t$ ,

$$\frac{dv_i(t)}{dt} = \left( \frac{dv_{i-1}(t)}{dt} - g_i(t) + g_{i+1}(t) \right)^+ \mathbb{I}_{\{g_i(t)=1\}}$$

for any  $i \geq 1$ . From this we see that, for a.e.  $t$ ,

$$\frac{dv_i(t)}{dt} = \begin{cases} \lambda & \text{if } i = 0 \\ \frac{dv_{i-1}(t)}{dt} & \text{if } i < m(\mathbf{g}(t)) \text{ and } i \geq 1, \\ \left( \frac{dv_{i-1}(t)}{dt} - 1 + g_{i+1}(t) \right)^+ & \text{if } i = m(\mathbf{g}(t)) \text{ and } i \geq 1, \\ 0 & \text{if } i > m(\mathbf{g}(t)). \end{cases} \quad (7.7)$$

and consequently  $p_j(\mathbf{g}(s)) = \frac{dv_j(s)}{ds} - \frac{dv_{j+1}(s)}{ds}$  for a.e.  $s$ . Substituting this back in (7.4) shows that  $\mathbf{g}$  solves the system of equations in (7.5). Conversely, for any solution  $\hat{\mathbf{g}}$  of (7.5), defining  $\hat{\mathbf{v}}$  by the right side of (7.7) by replacing  $\mathbf{g}$  with  $\hat{\mathbf{g}}$ , we see that  $(\hat{\mathbf{g}}, \hat{\mathbf{v}})$  solves (7.4). From the uniqueness result in Lemma 7 it then follows that in fact there is only one solution to the system of equations in (7.5) and this solution equals  $\mathbf{g}$  given in (7.3).

Consider now the time asymptotic behavior of  $\mathbf{g}$  given in (7.3). When  $\lambda < 1$ ,  $(\lambda, 0, 0 \dots) \in l_1$  is the unique fixed point of (7.3), as can be seen by setting the derivative of the right side of (7.5) to 0. In the critical case, i.e. when  $\lambda = 1$ , the situation is very different and in fact there are uncountably many fixed points given by  $\{\mathbf{f} \in l_1^\perp \mid m(\mathbf{f}) > 0, f_{m(\mathbf{f})+2} = 0\} \subset l_1^\perp$ , which once more is seen by checking that the derivative on the right side of (7.5) is 0 at exactly these points when  $\lambda = 1$ . In this chapter we are interested in the fluctuations of  $\mathbf{G}_n$  in the critical case when the system starts suitably close to one of the fixed points of (7.4). Thus for the remaining section we will assume that  $\lambda_n < 1$  for every  $n$  and  $\lambda_n \rightarrow 1$  as  $n \rightarrow \infty$ . In order to formulate precisely what we mean by ‘suitably close to the fixed point’ we need some definitions and notation. The functions  $\beta_n$  in the next definition will play a central role.

**Definition 8.** Define the function  $\beta_n : [0, 1] \rightarrow [0, 1]$  by

$$\beta_n(x) \doteq \prod_{i=0}^{d_n-1} \left( \frac{x - \frac{i}{n}}{1 - \frac{i}{n}} \right)^+ \quad (7.8)$$

The function  $\beta_n(\cdot)$  arises when sampling  $d_n$  random servers without replacement. Specifically, when  $nx \in \mathbb{N}$ ,  $\beta_n(x) = \mathbf{P}(\mathbb{A}_{n,d_n} \subseteq [nx]) = \binom{nx}{d_n} / \binom{n}{d_n}$ , where  $\mathbb{A}_{n,d_n}$  is a randomly chosen subset (without replacement) from  $[n]$  of size  $d_n$ . An alternative is to perform sampling with replacement, which corresponds to the simpler function  $\gamma_n(x) \doteq x^{d_n}$  in place of  $\beta_n$ .

We now introduce the notion of a ‘near fixed point’ of  $\mathbf{G}_n$ .

**Definition 9.** For  $n \in \mathbb{N}$ , the **near fixed point**  $\boldsymbol{\mu}_n$  of  $\mathbf{G}_n$  is the vector in  $l_1^\perp$  given as  $\boldsymbol{\mu}_n = (\mu_{n,1}, \mu_{n,2} \dots)$  where  $\mu_{n,i}$  are defined recursively as  $\mu_{n,1} = \lambda_n$  and  $\mu_{n,i+1} = \lambda_n \beta_n(\mu_{n,i})$  for  $i \geq 1$ .

Using  $\beta_n(x) \leq x^{d_n} \leq x$  and  $\lambda_n < 1$ , it is easy to check that  $\boldsymbol{\mu}_n \in l_1^\perp$ . The reason  $\boldsymbol{\mu}_n$  is referred to as a near fixed point of  $\mathbf{G}_n$  is discussed in Remark 7.3.1. To study the fluctuations of the process around the near fixed point  $\boldsymbol{\mu}_n$  we define the centered and scaled process,  $\mathbf{Z}_n$  as in (7.1). We now present our three main results on fluctuations which correspond to the three cases  $d_n/\sqrt{n} \rightarrow 0$ ,  $d_n/\sqrt{n} \rightarrow c \in (0, \infty)$ , and  $d_n/\sqrt{n} \rightarrow \infty$  respectively.

**Theorem 9.** Suppose that, as  $n \rightarrow \infty$ ,  $1 \ll d_n \ll \sqrt{n}$ ,  $\lambda_n \nearrow 1$ , and there is a  $k \in \mathbb{N}$  so that  $\mu_{n,k} \rightarrow 1$  and  $\beta'_n(\mu_{n,k}) \rightarrow \alpha \in [0, \infty)$  as  $n \rightarrow \infty$ . Further suppose that  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$  is tight and

that  $\mathbf{Z}_n(0) \xrightarrow{P} \mathbf{z}$  in  $l_2$ , where  $\mathbf{z}_{r+} = \mathbf{0}$  for some  $r > k$ . Then for any  $T \in (0, \infty)$ ,

$$\lim_{M \rightarrow \infty} \sup_n P\left(\|\mathbf{Z}_n\|_{2,T} > M\right) = 0. \quad (7.9)$$

Furthermore, if  $k > 1$ , then  $\sup_{t \in [\epsilon, T]} |Z_{n,i}(t)| \xrightarrow{P} 0$  as  $n \rightarrow \infty$  for any  $T < \infty$ ,  $0 < \epsilon \leq T$  and  $i \in [k-1]$ .

Consider the shifted process  $\mathbf{Y}_n(t) \doteq (\sum_{i=1}^k Z_{i,n}(t), Z_{k+1,n}(t), Z_{k+2,n}(t), \dots)$  and  $\mathbf{y} \doteq (\sum_{i=1}^k z_i, z_{k+1}, z_{k+2}, \dots)$ . Then  $\mathbf{Y}_n \Rightarrow \mathbf{Y}$  in  $\mathbb{D}([0, \infty) : l_2)$ , where  $\mathbf{Y} \in C([0, \infty) : l_2)$  is the unique pathwise solution to

$$\begin{aligned} Y_1(t) &= y_1 - (\alpha + \mathbb{I}_{\{k=1\}}) \int_0^t Y_1(s) ds + \int_0^t Y_2(s) ds + \sqrt{2}B(t) \\ Y_2(t) &= y_2 + \alpha \int_0^t Y_1(s) ds - \int_0^t Y_2(s) ds + \int_0^t Y_3(s) ds \\ Y_i(t) &= y_i - \int_0^t Y_i(s) ds + \int_0^t Y_{i+1}(s) ds \quad \text{for } i \in \{3, \dots, r-k+1\} \\ Y_i(t) &= 0 \quad \text{for } i > r-k+1, \end{aligned} \quad (7.10)$$

and  $B(\cdot)$  is a one dimensional standard Brownian motion.

**Remark 7.2.3.**

- (i) Note that the convergence  $\sup_{t \in [\epsilon, T]} |Z_{n,i}(t)| \xrightarrow{P} 0$  as  $n \rightarrow \infty$  for any  $0 < \epsilon \leq T$  is equivalent to the statement that  $Z_{n,i} \rightarrow 0$  in probability in  $\mathbb{D}((0, T] : \mathbb{R})$  where the latter space is equipped with the topology of uniform convergence on compacts. Note also that, since, for  $i \in [k-1]$ ,  $Z_{n,i}(0)$  may converge in general to a non-zero limit, the above convergence to 0 cannot be strengthened to a convergence in probability in  $\mathbb{D}([0, T] : \mathbb{R})$ .
- (ii) By Corollary 8 in Section 7.5, when  $\mu_{n,k}$  is away from 0,

$$\beta'_n(\mu_{n,k}) = (1 + o(1)) \frac{d_n \mu_{n,k+1}}{\lambda_n \mu_{n,k}}$$

as  $n \rightarrow \infty$ . Hence the assumptions  $d_n \rightarrow \infty$ ,  $\lambda_n \rightarrow 1$ ,  $\mu_{n,k} \rightarrow 1$  and  $\beta'_n(\mu_{n,k}) \rightarrow \alpha < \infty$  in Theorem 9 say that  $\mu_{n,k+1} \rightarrow 0$ . Since  $\mu_{n,k} \rightarrow 1$ , this in fact shows that  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_k$  in  $l_1^\perp$ , where recall that  $\mathbf{f}_k$  is one of the fixed points of the fluid-limit (7.3) when  $\lambda = 1$ . The fact that

the convergence happens in  $l_1^\downarrow$  can be seen on observing that if  $\mu_{n,k+1} \leq \epsilon$  then, by (7.37),  $\mu_{n,k+1+i} \leq \epsilon^{d_n^i}$ . This convergence, along with (7.9) shows that most queues will be of length  $k$  on any fixed interval  $[0, T]$ . We also note that in general  $\sqrt{n}(\boldsymbol{\mu}_n - \mathbf{f}_k)$  will diverge, and thus  $\sqrt{n}(\mathbf{G}_n - \mathbf{f}_k)$  will typically not be tight, in this regime.

- (iii) In the special case when the system starts sufficiently close to the near fixed point  $\boldsymbol{\mu}_n$  so that  $z_i = 0$  for  $i > k + 1$ , the limit process  $\mathbf{Y}$  simplifies to an essentially two dimensional process given as,  $Y_i(t) = 0$  for  $i > 2$ , and

$$\begin{aligned} Y_1(t) &= y_1 - (\alpha + \mathbb{I}_{\{k=1\}}) \int_0^t Y_1(s) ds + \int_0^t Y_2(s) ds + \sqrt{2}B(t) \\ Y_2(t) &= y_2 + \alpha \int_0^t Y_1(s) ds - \int_0^t Y_2(s) ds \end{aligned}$$

- (iv) The convergence behavior of  $\mathbf{Z}_n$  is governed by the sequence of parameters  $(d_n, \lambda_n)$ . In Corollary 9 from Section 7.5, we show that if  $1 \ll d_n^{k+1} \ll n$  and  $1 - \lambda_n = \frac{\xi_n + \log d_n}{d_n^k}$  with  $\xi_n \rightarrow -\log(\alpha) \in (-\infty, \infty]$  and  $\frac{\xi_n^2}{d_n} \rightarrow 0$ , then the conditions  $\mu_{n,k} \rightarrow 1$  and  $\beta'_n(\mu_{n,k}) \rightarrow \alpha \in [0, \infty)$  of Theorem 9 are satisfied. Using this fact we make the following observations. For simplicity, consider  $\mathbf{z} = 0$ .

- (a) Suppose that  $d_n = \log n$ ,  $1 - \lambda_n = \frac{\log \log n}{(\log n)^k}$ . In this case the assumptions of Theorem 9 are satisfied and one essentially sees non-zero fluctuations only in the  $k$ -th and  $k + 1$ -th coordinates. Note that as  $k$  becomes large, the traffic intensity increases and one sees more and more coordinates of the near fixed point approach 1.
- (b) With the same  $d_n$  as in (a) but a somewhat lower traffic intensity given as  $1 - \lambda_n = \frac{(\log n)^{1/2-\epsilon}}{(\log n)^k}$  for some  $\epsilon \in (0, 1/2)$ , one sees that condition of the theorem are satisfied with  $\alpha = 0$  (i.e.  $\beta'_n(\mu_{n,k}) \rightarrow 0$ ). Thus the limit process  $\mathbf{Y}$ , in the case  $k > 1$ , simplifies to  $Y_i = 0$  for  $i > 1$  and  $Y_1(t) = \sqrt{2}B(t)$ . When  $k = 1$ ,  $Z_1 = Y_1$  is instead given as the following Ornstein-Uhlenbeck(OU) process

$$Z_1(t) = - \int_0^t Z_1(s) ds + \sqrt{2}B(t). \quad (7.11)$$

- (c) With higher values of  $d_n$ , using Theorem 9, one can analyze fluctuations for systems with higher traffic intensity. For example, suppose that  $d_n = \frac{\sqrt{n}}{\log n}$ . Then the conditions of the theorem are satisfied with  $k = 1$  and  $1 - \lambda_n \sim (\log n)^2 / \sqrt{n}$ . In fact in this case  $\alpha = 0$  and the limit process is described by the one dimensional OU process (7.11). With a slightly higher traffic intensity given as  $1 - \lambda_n = ((\log n)^2 - 2 \log n \log \log n) / 2\sqrt{n}$  one obtains a two dimensional limit diffusion.
- (d) The theorem allows for traffic intensity in the Halfin-Whitt scaling regime (i.e.  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta > 0$ ) as well. Specifically, for  $k \geq 2$ , if  $d_n = (\sqrt{n} \log n)^{\frac{1}{k}}$  and  $(1 - \lambda_n) = \frac{\beta + o(1)}{\sqrt{n}}$  for some  $\beta > \beta_0 = 1/2k$ , the conditions of the theorem are satisfied with  $\alpha = 0$ . With slightly higher traffic intensity (e.g.  $\beta + o(1)$  replaced by  $\beta_0 + (\frac{1}{k} \log \log n - \log \alpha) / \log n$ ) conditions of the theorem are met with a non-zero  $\alpha$ .
- (e) Recall that a fixed point of (7.3) when  $\lambda = 1$  takes the form  $\mathbf{f}_k^\gamma \doteq \mathbf{f}_k + \gamma \mathbf{e}_k$ , where  $k \in \mathbb{N}$  and  $\gamma \in [0, 1)$ . Although Theorem 9 only considers settings where the near fixed point  $\boldsymbol{\mu}_n$  converges to  $\mathbf{f}_k^0 = \mathbf{f}_k$  for some  $k$ , it is possible to give conditions under which  $\boldsymbol{\mu}_n$  converges to a different fixed point. Specifically, suppose that  $1 \ll d_n^{k+1} \ll n$  and  $1 - \lambda_n = \frac{a}{d_n^k}$  for some  $a > 0$ . Then it can be checked using Lemma 22 that  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_k^\gamma$  with  $\gamma = e^{-a}$ .
- (v) Suppose for some  $a \in (0, \frac{1}{2})$ ,  $d_n = n^{a+o(1)}$  and  $\lambda_n$  is taken as in Remark 7.2.3 (iv) with  $k \in \mathbb{N}$  such that  $a(k+1) < 1$ . By Theorem 9, all but  $O(\sqrt{n})$  queues will have length  $k$  over bounded times. This result is analogous to (24, Theorem 1.1) which considers, for such choice of  $d_n, \lambda_n$ , the behavior of queues in equilibrium in a setting where  $d_n$  queues are sampled with replacement (instead of without replacement as in the current chapter). In fact, for this scenario, (24, Theorem 1.1) is able to show a stronger result which says that with high probability, as  $n \rightarrow \infty$ , most of the queues in equilibrium will have length  $k$  and that there will be no larger queues.

The next theorem describes the fluctuations of  $\mathbf{Z}_n$  when  $d_n$  is of order  $\sqrt{n}$ .

**Theorem 10.** *Suppose that  $\frac{d_n}{\sqrt{n}} \rightarrow c \in (0, \infty)$  and  $\lambda_n = 1 - \left( \frac{\log d_n}{d_n} + \frac{\alpha_n}{\sqrt{n}} \right)$  with  $\alpha_n \rightarrow \alpha \in (-\infty, \infty]$  and  $\alpha_n = o(n^{1/4})$ . Then,  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_1$  in  $l_1^\downarrow$ . Suppose further that  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$  is tight and  $\mathbf{Z}_n(0) \xrightarrow{P} \mathbf{z}$  in  $l_2$  with  $\mathbf{z}_{r+} = 0$  for some  $r \geq 2$ . Then, as  $n \rightarrow \infty$ ,  $\mathbf{Z}_n \Rightarrow \mathbf{Z}$  in  $\mathbb{D}([0, \infty) : l_2)$ , where  $\mathbf{Z}$  is the*

unique pathwise solution to:

$$\begin{aligned}
Z_1(t) &= z_1 - \int_0^t (Z_1(s) - Z_2(s))ds - (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds + \sqrt{2}B(t), \\
Z_2(t) &= z_2 - \int_0^t (Z_2(s) - Z_3(s))ds + (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds, \\
Z_i(t) &= z_i - \int_0^t (Z_i(s) - Z_{i+1}(s))ds \quad \text{for each } i \in \{3 \dots r\}, \\
Z_i(t) &= 0 \quad \text{for each } i > r,
\end{aligned}$$

and  $B$  is standard Brownian motion.

**Remark 7.2.4.**

- (i) Note that the coefficients in the above system of equations are only locally Lipschitz and have an exponential growth. However since  $c$  is positive, the system of equations has a unique pathwise solution as is shown in Lemma 39.
- (ii) Once more, when  $z_i = 0$  for all  $i > 2$ , the system of equations simplifies to a two dimensional system given as  $Z_i = 0$  for all  $i > 2$ , and

$$\begin{aligned}
Z_1(t) &= z_1 - \int_0^t (Z_1(s) - Z_2(s))ds - (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds + \sqrt{2}B(t), \\
Z_2(t) &= z_2 - \int_0^t Z_2(s)ds + (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds.
\end{aligned}$$

- (iii) In the regime considered in Theorem 10, the near fixed point  $\boldsymbol{\mu}_n$  can converge to only one particular fixed point of (7.3), namely  $\boldsymbol{f}_1$ . As before, the term  $\sqrt{n}(\boldsymbol{\mu}_n - \boldsymbol{f}_1)$  may diverge and thus  $\sqrt{n}(\boldsymbol{G}_n(\cdot) - \boldsymbol{f}_1)$  will in general not be tight.
- (iv) Suppose that  $d_n = c\sqrt{n}$  for some  $c > 0$ ,  $\boldsymbol{z} = 0$  and  $1 - \lambda_n = (\beta + o(1)) \log n / \sqrt{n}$  for some  $\beta > \beta_0 = 1/2c$ . Then the assumptions of the above theorem are satisfied with  $\alpha = \infty$  and the limit system simplifies to a one dimensional OU process given as  $Z_i = 0$  for all  $i > 1$ , and  $Z_1$  satisfies (7.11). If  $(\beta + o(1)) \log n$  is replaced by  $\beta_0 \log n + \gamma$  for some  $\gamma \in \mathbb{R}$ , we instead

obtain a two dimensional limit system given as  $Z_i = 0$  for all  $i > 2$ , and

$$\begin{aligned} Z_1(t) &= - \int_0^t (Z_1(s) - Z_2(s))ds - e^{-c\gamma} \int_0^t (e^{cZ_1(s)} - 1)ds + \sqrt{2}B(t), \\ Z_2(t) &= - \int_0^t Z_2(s)ds + e^{-c\gamma} \int_0^t (e^{cZ_1(s)} - 1)ds. \end{aligned}$$

Finally we consider the fluctuation behavior when  $d_n \gg \sqrt{n}$ . This time the limit system will involve reflected diffusion processes. Recall from (7.2) the definition of the Skorohod maps  $\Gamma_\alpha$  and  $\hat{\Gamma}_\alpha$  associated with a reflection barrier at  $\alpha \in \mathbb{R}$ . We will extend the definition of these maps to  $\alpha = \infty$  by setting

$$\Gamma_\infty(f) = f, \quad \hat{\Gamma}_\infty(f) = 0 \text{ for } f \in \mathbb{D}([0, \infty) : \mathbb{R}). \quad (7.12)$$

**Theorem 11.** *Suppose that  $\sqrt{n} \ll d_n$  and*

$$\lambda_n = 1 - \left( \frac{\log d_n}{d_n} + \frac{\alpha_n}{\sqrt{n}} \right), \text{ where } \alpha_n \rightarrow \alpha \in [0, \infty], \text{ with } \alpha_n^- = O(\sqrt{n}/d_n), \text{ and } \alpha_n = O(n^{1/6}). \quad (7.13)$$

*Then  $\mu_n \rightarrow \mathbf{f}_1$  in  $l_1^\downarrow$ . Suppose further that  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$  is tight and  $\mathbf{Z}_n(0) \xrightarrow{P} \mathbf{z}$  in  $l_2$  where  $z_1 \leq \alpha$  and  $\mathbf{z}_{r+} = 0$  for some  $r \geq 2$ . Then, as  $n \rightarrow \infty$ ,  $\mathbf{Z}_n \Rightarrow \mathbf{Z} \in \mathbb{D}([0, \infty) : l_2)$ , where  $(\mathbf{Z}, \eta)$  is a  $l_2 \times \mathbb{R}_+$  valued continuous process given as the unique solution to:*

$$\begin{aligned} Z_1(t) &= \Gamma_\alpha \left( z_1 - \int_0^t (Z_1(s) - Z_2(s))ds + \sqrt{2}B(\cdot) \right) (t), \\ Z_2(t) &= z_2 - \int_0^t (Z_2(s) - Z_3(s))ds + \eta(t), \\ \eta(t) &= \hat{\Gamma}_\alpha \left( z_1 - \int_0^t (Z_1(s) - Z_2(s))ds + \sqrt{2}B(\cdot) \right) (t), \\ Z_i(t) &= z_i - \int_0^t (Z_i(s) - Z_{i+1}(s))ds \quad \text{for each } i \in \{3 \dots r\}, \\ Z_i(t) &= 0 \quad \text{for each } i > r, \end{aligned} \quad (7.14)$$

*and  $B$  is a standard Brownian motion.*

As a corollary to this Theorem, we obtain the specific regime considered in (91) (in fact we provide a slight strengthening in that, unlike (91), we allow  $\alpha = 0$ ). See Remark 7.2.5 (v) for further discussion.



**Corollary 6.** *As  $n \rightarrow \infty$ , suppose that  $d_n \gg \sqrt{n} \log n$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \alpha \in [0, \infty)$ , along with  $\sqrt{n}(1 - \lambda_n) \geq (\sqrt{n} \log n)/d_n$  for large  $n$  if  $\alpha = 0$ . Let  $\mathbf{Y}_n(\cdot) \doteq \sqrt{n}(\mathbf{G}_n(\cdot) - \mathbf{f}_1)$  and assume that the sequence of random variables  $\{\|\mathbf{Y}_n(0)\|_1\}$  is tight, and as  $n \rightarrow \infty$ ,  $\mathbf{Y}_n(0) \xrightarrow{P} \mathbf{y} \in l_2$  with  $\mathbf{y}_{r+} = 0$  for some  $r \geq 2$ . Then  $\mathbf{Y}_n \Rightarrow \mathbf{Y}$  in  $\mathbb{D}([0, \infty) : l_2)$ , where  $(\mathbf{Y}, \tilde{\eta})$  is the  $l_2 \times [0, \infty)$  valued continuous process given by the unique solution to*

$$\begin{aligned} Y_1(t) &= \Gamma_0 \left( y_1 - \alpha \text{id}(\cdot) - \int_0^t (Y_1(s) - Y_2(s)) ds + \sqrt{2}B(\cdot) \right) (t) \\ Y_2(t) &= y_2 - \int_0^t (Y_2(s) - Y_3(s)) ds + \tilde{\eta}(t), \\ \tilde{\eta}(t) &= \hat{\Gamma}_0 \left( y_1 - \alpha \text{id}(\cdot) - \int_0^t (Y_1(s) - Y_2(s)) ds + \sqrt{2}B(\cdot) \right) (t), \\ Y_i(t) &= y_i - \int_0^t (Y_i(s) - Y_{i+1}(s)) ds \quad \text{for each } i \in \{3 \dots r\}, \\ Y_i(t) &= 0 \quad \text{for each } i > r, \end{aligned}$$

and  $B$  is a standard Brownian motion.

**Remark 7.2.5.**

- (i) The existence and uniqueness of solutions to the stochastic integral equations in (7.14) follows by standard fixed point arguments on using the Lipschitz property of the map  $\Gamma_\alpha$  on  $\mathbb{D}([0, \infty) : \mathbb{R})$ . This system of equations can equivalently be written as

$$\begin{aligned} Z_1(t) &= z_1 - \int_0^t (Z_1(s) - Z_2(s)) ds + \sqrt{2}B(t) - \eta(t), \\ Z_2(t) &= z_2 - \int_0^t (Z_2(s) - Z_3(s)) ds + \eta(t), \\ Z_i(t) &= z_i - \int_0^t (Z_i(s) - Z_{i+1}(s)) ds \quad \text{for each } i \in \{3 \dots r\}, \\ Z_i(t) &= 0 \quad \text{for each } i > r, \end{aligned} \tag{7.15}$$

where  $\eta = 0$  when  $\alpha = \infty$ , and when  $\alpha \in \mathbb{R}$ , it satisfies

$$\left. \begin{aligned} \eta(0) &= 0 \text{ and } \eta \text{ is a monotonically increasing function.} \\ Z_1(t) &\leq \alpha \\ \int_0^\infty (\alpha - Z_1(s)) d\eta(s) &= 0 \end{aligned} \right\} \tag{7.16}$$

- (ii) The convergence  $\mu_n \rightarrow \mathbf{f}_1$  along with tightness of  $\{\mathbf{Z}_n\}_{n \in \mathbb{N}}$  shows that, under the conditions of Theorems 10 or 11, most queues will be of length 1 on any fixed interval  $[0, T]$ .
- (iii) The limit system in Theorem 11 simplifies when  $z_i = 0$  for  $i > 2$  and is given as  $Z_i = 0$  for all  $i > 2$ , and

$$\begin{aligned} Z_1(t) &= z_1 - \int_0^t (Z_1(s) - Z_2(s))ds + \sqrt{2}B(t) - \eta(t), \\ Z_2(t) &= z_2 - \int_0^t Z_2(s)ds + \eta(t), \end{aligned}$$

where  $\eta$  is as in the statement of the theorem.

- (iv) Suppose that  $d_n = \sqrt{n} \log n / 2a$  for some  $a > 0$  and  $1 - \lambda_n = \frac{a}{\sqrt{n}} + \frac{2a(\log \log n + O(1))}{\sqrt{n} \log n}$ . Then the assumptions in Theorem 11 are satisfied with  $\alpha = 0$ . In this case the reflection barrier is at 0, namely  $Z_1(t) \leq 0$  for all  $t$ . Also note that since  $\sqrt{n}(1 - \lambda_n) \rightarrow a$ , we have that  $\mu_{n,1} = \lambda_n \rightarrow 1$ . Since  $d_n/\sqrt{n} \rightarrow \infty$ , this shows that for  $k \geq 2$

$$\sqrt{n}\mu_{n,2} = \sqrt{n}\lambda_n\beta_n(\lambda_n) \leq \sqrt{n}\lambda_n\lambda_n^{d_n} = \sqrt{n}(1 - (1 - \lambda_n))^{d_n+1} \rightarrow 0.$$

Using  $\mu_{n,i+1} \leq \mu_{n,i}^{d_n}$ , see that  $\sqrt{n}(\mu_n - \mathbf{f}_1) \rightarrow -a\mathbf{e}_1 \in l_1$  and hence the fluctuations of  $\mathbf{G}_n$  about the fixed point  $\mathbf{f}_1$  can be characterized as well. Specifically, letting  $\mathbf{Y}_n(\cdot) = \sqrt{n}(\mathbf{G}_n(\cdot) - \mathbf{f}_1) = \mathbf{Z}_n(\cdot) + \sqrt{n}(\mu_n - \mathbf{f}_1)$ , we see that, under the condition of the above theorem,  $\mathbf{Y}_n \Rightarrow \mathbf{Y}$  in  $\mathbb{D}([0, \infty) : l_2)$ , where  $\mathbf{Y} = \mathbf{Z} - a\mathbf{e}_1$  and hence, assuming  $z_i = 0$  for  $i > 2$ ,  $(\mathbf{Y}, \tilde{\eta}) \in C([0, \infty) : l_2 \times \mathbb{R}_+)$  is the unique solution to (7.16) with  $(Z_1, \eta, \alpha)$  replaced with  $(Y_1, \tilde{\eta}, -a)$ , and the equations

$$\begin{aligned} Y_1(t) &= y_1 - at - \int_0^t (Y_1(s) - Y_2(s))ds + \sqrt{2}B(t) - \tilde{\eta}(t), \\ Y_2(t) &= y_2 - \int_0^t Y_2(s)ds + \tilde{\eta}(t), \end{aligned}$$

where  $\mathbf{y} = \mathbf{z} - a\mathbf{e}_1$  and  $B$  is a standard Brownian motion. In particular, the limit  $\mathbf{Y}$  takes the same form as in (45, 91).

- (v) Suppose that  $d_n \gg \sqrt{n} \log n$ . Then it is easy to see that (7.13) holds with some  $\alpha > 0$  if and only if  $\sqrt{n}(1 - \lambda_n) \rightarrow \alpha > 0$ . This regime was studied in (91). Using the arguments as in (iv) above, it is easy to check that  $\sqrt{n}(\boldsymbol{\mu}_n - \mathbf{f}_1) \rightarrow -\alpha \mathbf{e}_1$  in  $l_1$  (and hence  $l_2$ ). Corollary 6 is immediate from this and Theorem 11. In particular we recover (91, Theorem 3). However the proof techniques in the current chapter are different from the stochastic coupling techniques employed in (91).
- (vi) Suppose  $\sqrt{n} \ll d_n \ll \sqrt{n} \log n$  and that (7.13) holds with  $\alpha < \infty$ . Then, as observed in (91), in this regime  $\mathbf{Y}_n$  is not tight. Indeed, it is easy to see that  $\sqrt{n}(1 - \lambda_n) = (\sqrt{n} \log d_n)/d_n + \alpha_n \rightarrow \infty$ . Nevertheless the process  $\mathbf{Z}_n$  converges in distribution and the limit process has a reflecting barrier at  $\alpha$ , i.e.  $Z_1 \leq \alpha$ . In particular, unlike the case  $d_n \gg \sqrt{n} \log n$ , the barrier in this case does not come from the constraint  $G_{n,1} \leq 1$ .
- (vii) Theorem 11 allows for a slower approach to criticality than  $n^{-1/2}$ , e.g.  $\lambda_n$  such that  $n^{1/3}(\lambda_n - 1) \rightarrow \gamma > 0$ . In this case  $\alpha = \infty$  and there is no reflection. When  $z_i = 0$  for all  $i \geq 1$ , this system reduces to the one dimensional OU process given by (7.11) with  $Z_i = 0$  for  $i > 1$ .

### 7.3 Poisson Representation of State Processes

We now embark on the proofs of the main results. We start with a brief overview of the organization of the proofs. In this Section we describe a specific construction of the state process. Proof of the law of large numbers (Theorem 8) is given in Section 7.4. Section 7.5 describes fine-scaled (deterministic) properties of the function  $\beta_n$  and the near fixed points  $\boldsymbol{\mu}_n$  which play a key technical role in the proofs of our diffusion approximations. Section 7.6 derives preliminary estimates required to prove all the main results for the fluctuations of the state process. Sections 7.7, 7.8 and 7.9 complete the proofs of Theorem 9, 10 and 11 respectively.

We start with a specific construction of the state process through time changed Poisson processes (cf. (66, 46)). A similar representation has been used in previous work on  $JSQ(d)$  systems (cf. (91, 45)). Let  $\{N_{i,+}, N_{i,-} : i \geq 1\}$  be a collection of mutually independent rate one Poisson processes given on some probability space  $(\Omega, \mathcal{F}, \mathbf{P})$ . Then  $\mathbf{G}_n$  has the following (equivalent in

distribution) representation. For  $i \geq 1$  and  $t \geq 0$

$$\begin{aligned} G_{n,i}(t) = G_{n,i}(0) &- \frac{1}{n} N_{i,-} \left( n \int_0^t [G_{n,i}(s) - G_{n,i+1}(s)] ds \right) \\ &+ \frac{1}{n} N_{i,+} \left( \lambda_n n \int_0^t [\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))] ds \right), \end{aligned} \quad (7.17)$$

where  $G_{n,0}(t) = 1$  for all  $t \geq 0$ . Denoting

$$A_{n,i}(t) \doteq N_{i,+} \left( \lambda_n n \int_0^t [\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))] ds \right), \quad D_{n,i}(t) \doteq N_{i,-} \left( n \int_0^t [G_{n,i}(s) - G_{n,i+1}(s)] ds \right),$$

the above evolution equation can be rewritten as

$$G_{n,i}(t) = G_{n,i}(0) - \frac{1}{n} D_{n,i}(t) + \frac{1}{n} A_{n,i}(t), \quad i \in \mathbb{N}, t \geq 0. \quad (7.18)$$

Here  $D_{n,i}$  describe events causing a decrease in  $G_{n,i}$  owing to completion of service events for jobs in queues of length exactly  $i$  whilst  $A_{n,i}$  describe events causing an increase in  $G_{n,i}$  which only occur if the chosen queue of a new job has exactly  $i - 1$  individuals; this occurs if amongst the  $d_n$  random choices made by this job, all of the chosen queues have load at least  $i - 1$  but not all have load at least  $i$ .

Let

$$\tilde{\mathcal{F}}_t^n = \sigma \{ A_i^n(s), D_i^n(s), s \leq t, i \geq 1 \},$$

and let  $\mathcal{F}_t^n$  be the augmentation of  $\tilde{\mathcal{F}}_t^n$  with  $\mathbf{P}$ -null sets. It then follows that, for each  $i \geq 1$

$$\begin{aligned} M_{n,i,+}(t) &\doteq \frac{1}{n} N_{i,+} \left( \lambda_n n \int_0^t \beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s)) ds \right) \\ &\quad - \lambda_n \int_0^t \beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s)) ds \end{aligned} \quad (7.19)$$

and

$$M_{n,i,-}(t) \doteq \frac{1}{n} N_{i,-} \left( n \int_0^t G_{n,i}(s) - G_{n,i+1}(s) ds \right) - \int_0^t (G_{n,i}(s) - G_{n,i+1}(s)) ds \quad (7.20)$$

are  $\{\mathcal{F}_t^n\}$ -martingales with predictable (cross) quadratic variation processes given, for  $t \geq 0$ , as

$$\begin{aligned}\langle M_{n,i,+} \rangle_t &= \frac{\lambda_n}{n} \int_0^t (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds, \quad i \geq 1, \\ \langle M_{n,i,-} \rangle_t &= \frac{1}{n} \int_0^t (G_{n,i}(s) - G_{n,i+1}(s)) ds, \quad i \geq 1, \\ \langle M_{n,i,-}, M_{n,j,-} \rangle_t &= 0, \quad \langle M_{n,i,+}, M_{n,j,+} \rangle_t = 0, \quad \text{for all } i, j \geq 1, i \neq j \text{ and} \\ \langle M_{n,i,+}, M_{n,k,-} \rangle_t &= 0 \text{ for all } i, k \geq 1.\end{aligned}$$

Using these martingales, the evolution of  $\mathbf{G}_n$  can be rewritten as

$$\begin{aligned}G_{n,i}(t) &= G_{n,i}(0) - \int_0^t (G_{n,i}(s) - G_{n,i+1}(s)) ds \\ &\quad + \lambda_n \int_0^t (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds + M_{n,i}(t), \quad i \geq 1\end{aligned}\tag{7.21}$$

where  $M_{n,i}(t) \doteq M_{n,i,+}(t) - M_{n,i,-}(t)$  and

$$\langle M_{n,i} \rangle_t = \frac{1}{n} \left( \int_0^t (G_{n,i}(s) - G_{n,i+1}(s)) ds + \lambda_n \int_0^t (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds \right). \tag{7.22}$$

We will assume throughout that  $\mathbf{G}_n(0) \in l_1^\downarrow$  a.s. Then it follows that, for every  $t \geq 0$ ,  $\|\mathbf{G}_n(t)\|_1 < \infty$  almost surely. Indeed, over any time interval  $[0, t]$  finitely many jobs enter the system a.s. and denoting by  $k(n)$  the number of jobs that arrive over  $[0, t]$ , we see that  $\|\mathbf{G}_n(t)\|_1 \leq \|\mathbf{G}_n(0)\|_1 + k(n)/n < \infty$  a.s. Thus  $\mathbf{G}_n$  is a stochastic process with sample paths in  $\mathbb{D}([0, \infty) : l_1^\downarrow)$ . Note that, for any  $t > 0$ ,  $\|\mathbf{G}_n(t) - \mathbf{G}_n(t-)\|_1 \leq \frac{1}{n}$ .

**Remark 7.3.1.** Let  $\mathbf{a}_n, \mathbf{b} : l_1^\downarrow \rightarrow l_1$  be given by

$$\mathbf{a}_n(\mathbf{x})_i \doteq \lambda_n(\beta_n(x_{i-1}) - \beta_n(x_i)), \quad \mathbf{b}(\mathbf{x})_i \doteq x_i - x_{i+1}, \quad \mathbf{x} \in l_1^\downarrow, \quad i \geq 1,$$

where, by convention, for  $\mathbf{x} \in l_1^\downarrow$ ,  $x_0 = 1$ . Then (7.21) can be rewritten as an evolution equation in  $l_1$  as,

$$\mathbf{G}_n(t) = \mathbf{G}_n(0) + \int_0^t [\mathbf{a}_n(\mathbf{G}_n(s)) - \mathbf{b}(\mathbf{G}_n(s))] ds + \mathbf{M}_n(t), \tag{7.23}$$

where  $\mathbf{M}_n(t) \doteq (M_{n,i}(t))_{i \geq 1}$  is a stochastic process with sample paths in  $\mathbb{D}([0, \infty) : l_1)$  and the integral is a Bochner-integral (133). Note that the near fixed point  $\boldsymbol{\mu}_n$  from Definition 9 satisfies  $\mathbf{a}_n(\boldsymbol{\mu}_n) = \mathbf{b}(\boldsymbol{\mu}_n)$ . It is in fact the unique solution to,

$$\mathbf{a}_n(\mathbf{x}) = \mathbf{b}(\mathbf{x}) \quad \text{for } \mathbf{x} \in l_1^\downarrow, \quad (7.24)$$

as is seen by adding up all the coordinates of (7.24) and using  $\mathbf{x} \in l_1$ . In Lemma 15 we will see that for any  $T > 0$ , as  $n \rightarrow \infty$ ,  $\sup_{t \leq T} \|\mathbf{M}_n(t)\|_2 \xrightarrow{P} 0$ . Hence if  $\mathbf{G}_n(0) = \boldsymbol{\mu}_n$ , then by (7.23), we expect the process  $\mathbf{G}_n(t)$  to stay close to  $\boldsymbol{\mu}_n$  (over any compact time interval) as  $n \rightarrow \infty$ . In this sense  $\boldsymbol{\mu}_n$  can be viewed as a ‘near fixed point’ of  $\mathbf{G}_n(\cdot)$  and the terminology in Definition 9 is justified. Another reason for this terminology comes from the results in Theorems 9–11 which show that, under conditions,  $\boldsymbol{\mu}_n$  converges to one of the fixed points of the fluid limit (7.3) when  $\lambda = 1$ .

## 7.4 The Law of Large Numbers

In this section we prove Proposition 7 and Theorem 8.

### 7.4.1 Uniqueness of Fluid Limit Equations.

In this subsection we show that there is at most one solution of (7.3) in  $C([0, \infty) : l_1^\downarrow \times l_\infty)$ . Results of Section 7.4.2 will provide existence of solutions to this equation. Suppose  $(\mathbf{g}, \mathbf{v})$  and  $(\mathbf{g}', \mathbf{v}')$  are two solutions to (7.3) in  $C([0, \infty) : l_1^\downarrow \times l_\infty)$ . We will now argue that the two solutions are equal.

We claim that that  $v'_i$  and  $v_i$  are non-zero for only finitely many  $i$ ’s. Indeed, since  $\mathbf{g}, \mathbf{g}' \in C([0, T] : l_1^\downarrow)$ , there is a constant  $C \in (0, \infty)$  so that  $\sup_{s \leq T} \|\mathbf{g}(s)\|_1 \vee \sup_{s \leq T} \|\mathbf{g}'(s)\|_1 \leq C$ . Since

$$x_i \leq \|\mathbf{x}\|_1 / i \text{ for any } \mathbf{x} \in l_1^\downarrow, \quad (7.25)$$

taking  $M \doteq \lceil C + 1 \rceil \in \mathbb{N}$  shows that  $\sup_{s \leq T} g_i(s) \vee g'_i(s) < 1$  for any  $i \geq M$ . But then by the equivalent representation of (7.3) given in (7.4) (in particular the second line), we must have  $v_i = v'_i = 0$  for any  $i \geq M$ . This proves the claim.

Since  $v_i = v'_i = 0$  for  $i \geq M$ , the first line of the equivalent formulation in (7.4) shows that both  $\mathbf{x} = \mathbf{g}$  and  $\mathbf{x} = \mathbf{g}'$  satisfy the integral equations

$$x_i(t) = r_i - \int_0^t (x_i(s) - x_{i+1}(s)) ds \quad \text{for } i \geq M+1 \text{ and } t \in [0, T].$$

By standard arguments using Gronwall's lemma (46, Appendix 5), we then must have  $g_i = g'_i$  for each  $i \geq M+1$ . Indeed, letting  $z_i(\cdot) \doteq g_i(\cdot) - g'_i(\cdot)$  for  $i \geq M+1$  and  $v(t) \doteq \sum_{i=M+1}^{\infty} |z_i(t)|$  for  $t \in [0, T]$ , we have that

$$|z_i(t)| \leq \int_0^t (|z_i(s)| + |z_{i+1}(s)|) ds \quad \text{for all } i \geq M+1, \text{ and } t \in [0, T]$$

and so

$$v(t) \leq 2 \int_0^t v(s) ds, \quad t \in [0, T],$$

which implies that  $v(t) = 0$  for  $t \in [0, T]$ .

We now show that  $g_i = g'_i$  for  $i \leq M$ . From the definition of the Skorohod map in (7.2) we see that for  $f_1, f_2 \in \mathbb{D}([0, \infty) : \mathbb{R})$  with  $f_i(0) \leq 1$ ,  $i = 1, 2$ , and  $t \geq 0$

$$\|\Gamma_1(f_1) - \Gamma_1(f_2)\|_{*,t} \leq 2 \|f_1 - f_2\|_{*,t}, \quad \|\hat{\Gamma}_1(f_1) - \hat{\Gamma}_1(f_2)\|_{*,t} \leq \|f_1 - f_2\|_{*,t}.$$

Thus, since  $(\mathbf{g}, \mathbf{v})$  and  $(\mathbf{g}', \mathbf{v}')$  solve (7.3),

$$\|g_i - g'_i\|_{*,t} \leq 2 \left( \int_0^t \|g_i - g'_i\|_{*,s} ds + \int_0^t \|g_{i+1} - g'_{i+1}\|_{*,s} ds + \|v_{i-1} - v'_{i-1}\|_{*,t} \right), \quad \text{and} \quad (7.26)$$

$$\|v_i - v'_i\|_{*,t} \leq \int_0^t \|g_i - g'_i\|_{*,s} ds + \int_0^t \|g_{i+1} - g'_{i+1}\|_{*,s} ds + \|v_{i-1} - v'_{i-1}\|_{*,t} \quad (7.27)$$

for any  $i \geq 1$ . Let  $H_t \doteq \max_{i \in \{1, \dots, M\}} \|g_i - g'_i\|_{*,t}$ . Note  $g_{M+1} = g'_{M+1}$  and hence  $H_t = \max_{i \in \{1, \dots, M+1\}} \|g_i - g'_i\|_{*,t}$ . Then from (7.27), we have

$$\|v_i - v'_i\|_{*,t} \leq 2 \int_0^t H_s ds + \|v_{i-1} - v'_{i-1}\|_{*,t} \quad \text{for any } i \leq M. \quad (7.28)$$

Repeatedly using (7.28) along with  $v_0 = v'_0$  shows that  $\|v_i - v'_i\|_{*,t} \leq 2i \int_0^t H_s ds$  for any  $i \leq M$ .

Using this bound in (7.26) shows for  $1 \leq i \leq M$ :

$$\|g_i - g'_i\|_{*,t} \leq 2 \left( 2 \int_0^t H_s ds + 2(i-1) \int_0^t H_s ds \right) = 4i \int_0^t H_s ds.$$

Hence considering the maximum of  $\|g_i - g'_i\|_{*,t}$  over  $1 \leq i \leq M$  we get

$$0 \leq H_t \leq 4M \int_0^t H_s ds \quad \text{for each } t \in [0, T].$$

Gronwall's Lemma now shows that  $H_T = 0$ , and hence  $g_i = g'_i$  for  $i = 1 \dots M$ . Finally, since  $v_0 = v'_0$ , we see recursively from the second equation in (7.3) that  $v_i = v'_i$  for all  $i \geq 0$ .  $\square$

#### 7.4.2 Tightness and Limit Point Characterization

Some of the arguments in this section are similar to (91) however in order to keep the presentation self-contained we provide details in a concise manner. The next result establishes the convergence of the martingale term  $\mathbf{M}_n$  in the semimartingale decomposition in (7.23). Throughout this subsection and the next we assume that the conditions of Theorem 8 are satisfied, namely,  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{r}$  in  $l_1^\downarrow$ ,  $\lambda_n \rightarrow \lambda$  and  $d_n \rightarrow \infty$ , as  $n \rightarrow \infty$ .

**Lemma 15.** *For any  $T > 0$ ,  $\sup_{s \leq T} \|\mathbf{M}_n(s)\|_2 \xrightarrow{P} 0$ .*

*Proof.* It suffices to show that for any  $T > 0$ ,  $\lim_n \mathbf{E} \sup_{s \leq T} \|\mathbf{M}_n(s)\|_2^2 = 0$ . Applying Doob's maximal inequality we have that

$$\mathbf{E} \sup_{s \leq T} \|\mathbf{M}_n(s)\|_2^2 \leq 4\mathbf{E} \|\mathbf{M}_n(T)\|_2^2 = 4\mathbf{E} \sum_{i \geq 1} M_{n,i}(T)^2. \quad (7.29)$$

Since  $\mathbf{E} M_{n,i}^2(T) = \mathbf{E} \langle M_{n,i} \rangle_T$ , using the monotone convergence theorem in (7.29) shows,

$$\mathbf{E} \sup_{s \leq T} \|\mathbf{M}_n(s)\|_2^2 \leq 4\mathbf{E} \sum_{i \geq 1} \langle M_{n,i} \rangle_T \leq 4 \frac{T(1 + \sup_n \lambda_n)}{n}, \quad (7.30)$$



where the last inequality is from (7.22) on observing that

$$\sum_{i=1}^{\infty} \langle M_{n,i} \rangle_T \leq \frac{1}{n} \int_0^T G_{n,1}(s) + \frac{\lambda_n}{n} \int_0^T \beta_n(G_{n,0}(t)) \leq \frac{T(1 + \lambda_n)}{n}.$$

Sending  $n \rightarrow \infty$  in (7.30) completes the proof of the lemma.  $\square$

The next lemma characterizes compact sets in  $l_1^\downarrow$ . The proof is standard and can be found for example in (91).

**Proposition 12.** A subset  $C \subseteq l_1^\downarrow$  is precompact if and only if the following two conditions hold:

1. (norm-bounded)  $\sup_{\mathbf{x} \in C} \|\mathbf{x}\|_1 < \infty$ , and
2. (uniformly decaying tails)  $\limsup_{M \rightarrow \infty} \sup_{\mathbf{x} \in C} \sum_{i > M} |x_i| = 0$ .

**Lemma 16.** For each  $n \in \mathbb{N}$  there is a square integrable  $\{\mathcal{F}_t^n\}$ -martingale  $\{L_n(t)\}$  such that, for any  $t \geq 0$ ,

$$\sup_{s \in [0, t]} \|\mathbf{G}_n(s)\|_1 \leq \|\mathbf{G}_n(0)\|_1 + \lambda_n t + L_n(t).$$

Furthermore,  $\langle L_n \rangle_t \leq \frac{\lambda_n t}{n}$ , for all  $t \geq 0$ .

*Proof.* For  $i = 1, \dots, n$ , let  $X_i(t)$  denote the number of jobs in the  $i$ -th server's queue at time  $t$ .

Then

$$\|\mathbf{G}_n(t)\|_1 = \sum_{j=1}^{\infty} G_{n,j}(t) = \sum_{j=1}^{\infty} \sum_{i=1}^n \frac{\mathbb{I}_{\{X_i(t) \geq j\}}}{n} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{\infty} \mathbb{I}_{\{X_i(t) \geq j\}} = \frac{1}{n} \sum_{i=1}^n X_i(t).$$

Hence  $\|\mathbf{G}_n(t)\|_1$  is the total number of jobs in the system at time  $t$ , divided by  $n$ .

Since the total number of jobs in the system at time  $t$  is bounded above by the sum of number of job arrivals by time  $t$  and the initial number of jobs,  $\sup_{s \in [0, t]} \|\mathbf{G}_n(s)\|_1 \leq \|\mathbf{G}_n(0)\|_1 + \frac{A_n(t)}{n}$ , where  $A_n(t)$  is the total number of arrivals to the system by time  $t$ . Since,  $A_n$  is a Poisson process with arrival rate  $\lambda_n n$ , the result follows on setting  $L_n(t) = \frac{A_n(t)}{n} - \lambda_n t$ ,  $t \geq 0$ .  $\square$

The estimate in the next lemma will be useful when applying Aldous-Kurtz tightness criteria (46) for proving tightness of  $\{\mathbf{G}_n\}$ .

**Lemma 17.** Fix  $n \in \mathbb{N}$  and  $\delta \in (0, \infty)$ . Let  $\tau$  be a bounded  $\{\mathcal{F}_t^n\}$ -stopping time. Then

$$\mathbf{E} \|\mathbf{G}_n(\tau + \delta) - \mathbf{G}_n(\tau)\|_1 \leq (\lambda_n + 1)\delta$$

*Proof.* From (7.18), for any  $i \in \mathbb{N}$ ,

$$|G_{n,i}(\tau + \delta) - G_{n,i}(\tau)| \leq \frac{1}{n}(A_{n,i}(\tau + \delta) - A_{n,i}(\tau) + D_{n,i}(\tau + \delta) - D_{n,i}(\tau)). \quad (7.31)$$

From (7.19) and (7.20) we see that

$$\begin{aligned} \mathbf{E} \frac{1}{n}(A_{n,i}(\tau + \delta) - A_{n,i}(\tau)) &= \lambda_n \mathbf{E} \int_{\tau}^{\tau+\delta} (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds \\ \mathbf{E} \frac{1}{n}(D_{n,i}(\tau + \delta) - D_{n,i}(\tau)) &= \mathbf{E} \int_{\tau}^{\tau+\delta} (G_{n,i}(s) - G_{n,i+1}(s)) ds. \end{aligned}$$

Using the above identities in (7.31)

$$\begin{aligned} \mathbf{E} |G_{n,i}(\tau + \delta) - G_{n,i}(\tau)| \\ \leq \lambda_n \mathbf{E} \int_{\tau}^{\tau+\delta} (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds + \mathbf{E} \int_{\tau}^{\tau+\delta} (G_{n,i}(s) - G_{n,i+1}(s)) ds \end{aligned} \quad (7.32)$$

Adding (7.32) over various values of  $i \in \mathbb{N}$ , we have

$$\begin{aligned} \mathbf{E} \|\mathbf{G}_n(\tau + \delta) - \mathbf{G}_n(\tau)\|_1 &\leq \lambda_n \sum_{i=1}^{\infty} \mathbf{E} \int_{\tau}^{\tau+\delta} (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds \\ &\quad + \sum_{i=1}^{\infty} \mathbf{E} \int_{\tau}^{\tau+\delta} (G_{n,i}(s) - G_{n,i+1}(s)) ds \\ &\leq \mathbf{E} \int_{\tau}^{\tau+\delta} (\lambda_n \beta_n(G_{n,0}(s)) + G_{n,1}(s)) ds \\ &\leq (\lambda_n + 1)\delta. \end{aligned}$$

□

The following lemma will be useful in verifying the tightness of  $\{\mathbf{G}_n(t)\}$  in  $l_1^1$  for each fixed  $t \geq 0$ .

**Lemma 18.** For every  $n, m \in \mathbb{N}$  there is a square integrable  $\{\mathcal{F}_t^n\}$  martingale  $L_{n,m}(\cdot)$  so that, for all  $t \geq 0$ ,

$$\sup_{s \leq t} \sum_{i > m} G_{n,i}(s) \leq \sum_{i > m} G_{n,i}(0) + \frac{\lambda_n t}{m} \|\mathbf{G}_n\|_{1,t} + L_{n,m}(t)$$

and  $\langle L_{n,m} \rangle_t \leq \frac{\lambda_n t}{nm} \|\mathbf{G}_n\|_{1,t}$ .

*Proof.* From (7.17), for any  $i \in \mathbb{N}$  and  $t \geq 0$ :

$$G_{n,i}(t) \leq G_{n,i}(0) + \frac{1}{n} N_{+,i} \left( n \lambda_n \int_0^t \beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s)) ds \right) \quad (7.33)$$

Consider the point-process given by

$$B_{n,m}(t) \doteq \sum_{i > m} N_{+,i} \left( n \lambda_n \int_0^t \beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s)) ds \right).$$

Adding over  $i > m$  in (7.33) we get

$$\sup_{s \leq t} \sum_{i > m} G_{n,i}(s) \leq \sum_{i > m} G_{n,i}(0) + \frac{1}{n} B_{n,m}(t) \quad (7.34)$$

It is easy to see that, with

$$b_{n,m}(t) \doteq n \lambda_n \sum_{i > m} \int_0^t \beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s)) ds, \quad t \geq 0,$$

$\tilde{L}_{n,m}(t) \doteq B_{n,m}(t) - b_{n,m}(t)$  is a  $\mathcal{F}_t^n$ -martingale and

$$\begin{aligned} \langle \tilde{L}_{n,m} \rangle_t &= b_{n,m}(t) = n \lambda_n \int_0^t \beta_n(G_{n,m}(s)) ds \\ &\leq n \lambda_n \int_0^t G_{n,m}(s) ds \leq n \lambda_n t \left( \sup_{s \leq t} G_{n,m}(s) \right) \leq \frac{n \lambda_n t}{m} \|\mathbf{G}_n\|_{1,t}, \end{aligned}$$

where, for the last inequality we have used (7.25). The lemma now follows on setting  $L_{n,m}(t) = \tilde{L}_{n,m}(t)/n$  and using (7.34).  $\square$

Recall that under our assumptions,  $\lambda_n \rightarrow \lambda$  and  $d_n \rightarrow \infty$  as  $n \rightarrow \infty$ .

**Lemma 19.** *Suppose that  $\{\mathbf{G}_n(0)\}_{n \geq 1}$  is a tight sequence of  $l_1^\downarrow$  valued random variables. Then for any  $T > 0$ ,  $\{\mathbf{G}_n\}_{n \geq 1}$  is a tight sequence of  $\mathbb{D}([0, T] : l_1^\downarrow)$  valued random variables.*

*Proof.* To show that  $\{\mathbf{G}_n\}_{n \geq 1}$  is tight it suffices to show that (cf. (67, Theorem 2.7))

1. For any  $t \in [0, T]$  and  $\epsilon > 0$ , there is a compact set  $\Gamma \subseteq l_1^\downarrow$  so that  $\inf_{n \in \mathbb{N}} \mathbf{P}(\mathbf{G}_n(t) \in \Gamma) \geq 1 - \epsilon$ .
2.  $\lim_{\delta \rightarrow 0} \limsup_{n \rightarrow \infty} \sup_{\tau \leq T} \mathbf{E} \|\mathbf{G}_n(\tau + \delta) - \mathbf{G}_n(\tau)\|_1 = 0$ , where the innermost supremum is taken over all  $\mathcal{F}_t^n$ -stopping times  $\tau$  that are bounded by  $T - \delta$ .

The second condition is immediate from Proposition 17. Now consider (1). Fix  $\epsilon > 0$ . Let  $\bar{\lambda} = \sup_{n \geq 1} \lambda_n$ . Since  $\mathbf{G}_n(0)$  is tight, there is a compact  $K_1 \subset l_1^\downarrow$  such that

$$\mathbf{P}(\mathbf{G}_n(0) \in K_1) \geq 1 - \frac{\epsilon}{8} \text{ for all } n \in \mathbb{N}.$$

From Proposition 12 there is a  $\kappa_1 \in (0, \infty)$  such that  $\sup_{\mathbf{x} \in K_1} \|\mathbf{x}\|_1 \leq \kappa_1$ . From Lemma 16 we can find  $\kappa_2 \in (0, \infty)$  so that

$$\mathbf{P}(\bar{\lambda}T + \|L_n\|_{1,T} > \kappa_2) \leq \frac{\epsilon}{8}.$$

Then, using the above estimates and Lemma 16 again, with  $\kappa = \kappa_1 + \kappa_2$ ,

$$\mathbf{P}(\|\mathbf{G}_n\|_{1,T} \geq \kappa) \leq \frac{\epsilon}{4}.$$

Let  $m_k \uparrow \infty$  be a sequence such that  $4 \frac{\bar{\lambda}T\kappa}{m_k^{1/2}} \leq \frac{\epsilon}{2^{k+2}}$  for all  $k \in \mathbb{N}$ . Define

$$K_2 = \left\{ \mathbf{y} \in l_1^\downarrow : \|\mathbf{y}\|_1 \leq \kappa \text{ and for some } \mathbf{x} \in K_1, \sum_{i > m_k} y_i \leq \sum_{i > m_k} x_i + \frac{\bar{\lambda}T\kappa}{m_k} + \frac{1}{m_k^{1/4}}, \forall k \in \mathbb{N} \right\}.$$

Since  $K_1$  is compact, it is immediate from Proposition 12 that  $K_2$  is precompact in  $l_1^\downarrow$ . Also, using Lemma 18, for any  $t \in [0, T]$ ,

$$\begin{aligned} \mathbf{P}(\mathbf{G}_n(t) \in K_2^c) &\leq \mathbf{P}(\|\mathbf{G}_n\|_{1,T} \geq \kappa) + \mathbf{P}(\mathbf{G}_n(0) \in K_1^c) + \mathbf{P}(\|L_{n,m_k}\|_{*,T} > \frac{1}{m_k^{1/4}} \text{ for some } k \in \mathbb{N}) \\ &\leq \frac{\epsilon}{4} + \frac{\epsilon}{8} + 4\kappa\bar{\lambda}T \sum_{k=1}^{\infty} m_k^{1/2} \frac{1}{m_k} \leq \epsilon, \end{aligned}$$

where the second inequality follows from Doob's maximal inequality and from the expression of  $\langle L_{n,m_k} \rangle$  in Lemma 18 and the third inequality follows from the choice of  $\{m_k\}$ . This proves (1) and completes the proof of the lemma.  $\square$

The following lemma gives a characterization of the limit points of  $\mathbf{G}_n$ .

**Lemma 20.** *Fix  $T \in (0, \infty)$ . Suppose that, along some subsequence  $\{n_k\}_{k \geq 1}$ ,  $\mathbf{G}_{n_k} \Rightarrow \mathbf{G}$  in  $\mathbb{D}([0, T] : l_1^\downarrow)$  as  $k \rightarrow \infty$ . Then  $\mathbf{G} \in C([0, T] : l_1^\downarrow)$  a.s., and (7.3) is satisfied with  $(g_i, v_i)$  replaced with  $(G_i, V_i)$ , where  $V_i$  are defined recursively using the second equation in (7.3) with  $V_0(t) = \lambda t$  for  $t \geq 0$ .*

*Proof.* From Lemma 15 we see that  $\mathbf{M}_{n_k} \xrightarrow{P} 0$ , in  $\mathbb{D}([0, T] : l_2)$ . By Skorohod embedding theorem, let us assume that  $\mathbf{G}_{n_k}, \mathbf{M}_{n_k}, \mathbf{G}$  are all defined on the same probability space and

$$(\mathbf{G}_{n_k}, \mathbf{M}_{n_k}) \rightarrow (\mathbf{G}, 0), \text{ a.s.}$$

in  $\mathbb{D}([0, T] : l_1^\downarrow \times l_2)$ . Since the jumps of  $\mathbf{G}_n$  have size at most  $1/n$ ,  $\mathbf{G}$  is continuous and  $\|\mathbf{G}(s) - \mathbf{G}_{n_k}(s)\|_{1,T} \rightarrow 0$  a.s. Similarly,  $\|\mathbf{M}_{n_k}(s)\|_{2,T} \rightarrow 0$  almost surely. To simplify notation from now on we will take  $n_k = n$ .

Let  $V_{n,i}(t) \doteq \lambda_n \int_0^t \beta_n(G_{n,i}(s)) ds$  for  $i \geq 1$  and  $V_{n,0}(t) \doteq \lambda_n t$ . From (7.17), for any  $i \geq 1$

$$G_{n,i}(t) = G_{n,i}(0) - \int_0^t (G_{n,i}(s) - G_{n,i+1}(s)) ds + V_{n,i-1}(t) - V_{n,i}(t) + M_{n,i}(t). \quad (7.35)$$

For  $i \in \mathbb{N}$ ,  $\sup_{s \leq T} |G_{n,i}(s) - G_i(s)| \leq \sup_{s \leq T} \|\mathbf{G}_n(s) - \mathbf{G}(s)\|_1 \rightarrow 0$  and  $\sup_{s \leq T} |M_{n,i}(s)| \leq \sup_{s \leq T} \|\mathbf{M}_n(s)\|_2 \rightarrow 0$ , almost surely as  $n \rightarrow \infty$ . We now show that, for each  $i \in \mathbb{N}_0$ ,  $V_{n,i}$  converges uniformly (a.s.) to some limit process  $V_i$ . Clearly this is true for  $i = 0$  and in fact  $V_0(t) = \lambda t$ ,  $t \geq 0$ . Proceeding recursively, suppose now that  $V_{n,i-1} \rightarrow V_{i-1}$  for some  $i \geq 1$ . Then, since all the terms in (7.35), except  $V_{n,i}$ , converge uniformly,  $V_{n,i}$  must converge uniformly as well to some limit process  $V_i$ . Sending  $n \rightarrow \infty$  in (7.35) we get, for every  $t \leq T$  and  $i \geq 1$ :

$$G_i(t) = G_i(0) - \int_0^t (G_i(s) - G_{i+1}(s)) ds + V_{i-1}(t) - V_i(t), \text{ a.s.}$$

This shows the first line in (7.4) is satisfied with  $(g_i, v_i)$  replaced with  $(G_i, V_i)$ .

We now show that the second line in (7.4) is satisfied as well. Since  $V_i$  is the limit of  $\{V_{n,i}\}$ , the following properties hold:

- (i)  $V_0(t) = \lambda t$  for all  $t \in [0, T]$ .
- (ii)  $V_i$  is continuous, non-decreasing and  $V_i(0) = 0$ .
- (iii) For any  $t \in [0, T]$ ,  $\int_0^t (1 - G_i(s)) dV_i(s) = 0$ . This is a consequence of the following identities:

$$\begin{aligned}
\int_0^t (1 - G_i(s)) dV_i(s) &= \lim_n \int_0^t (1 - G_i(s)) dV_{n,i}(s) \\
&= \lim_n \int_0^t \lambda_n (1 - G_i(s)) \beta_n(G_{n,i}(s)) ds \\
&= \lambda \int_0^t \lim_{n \rightarrow \infty} (1 - G_i(s)) \beta_n(G_{n,i}(s)) ds \\
&= 0
\end{aligned}$$

where the first equality holds since  $G_i$  is a continuous and bounded function and  $V_{n,i} \rightarrow V_i$  uniformly on  $[0, T]$ ; the second equality uses the definition of  $V_{n,i}$ , the third is from the dominated convergence theorem, and the fourth follows since  $\beta_n(x) \leq x^{d_n}$ , for  $x \in [0, 1]$  and  $d_n \rightarrow \infty$ ,  $\beta_n(x) \rightarrow 0$  for every  $x \in [0, 1]$ .

Thus we have verified that the second line in (7.4) is satisfied with  $(G_i, V_i)$  as well. The result is now immediate from Remark 7.2.1.  $\square$

### 7.4.3 Completing the Proof of LLN

We can now complete the proofs of Proposition 7 and Theorem 8.

*Proof of Proposition 7.* Fix  $\mathbf{r} \in l_1^\downarrow$ ,  $\lambda > 0$  and choose a sequence  $\mathbf{r}_n \in l_1^\downarrow$  such that  $\mathbf{r}_n \rightarrow \mathbf{r}$  in  $l_1^\downarrow$  and for each  $i$ ,  $nr_{n,i} \in \mathbb{N}_0$ . Consider parameters  $\lambda_n = \lambda$ ,  $d_n = n$  and a JSQ( $d_n$ ) system initialized at  $\mathbf{G}_n(0) = \mathbf{r}_n$ . From Lemma 20 we have that there is at least one solution of (7.3) which is given as a limit point of an arbitrary weakly convergent subsequence of  $\mathbf{G}_n$  (such a sequence exists in view of the tightness shown in Lemma 19). The fact that this equation can have at most one solution was shown in Section 7.4.1. The result follows.  $\square$

*Proof of Theorem 8.* Since  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{r}$  in  $l_1^\downarrow$ , the hypothesis of Lemma 19 is satisfied, and thus the sequence  $\{\mathbf{G}_n\}_{n \geq 1}$  is tight in  $\mathbb{D}([0, T] : l_1^\downarrow)$  for any fixed  $T > 0$ . The result is now immediate from Lemma 20 and unique solvability of (7.3) shown in Proposition 7.  $\square$

**Remark 7.4.1.** We note that the proofs of Lemma 20 and Theorem 8 also show that, under the conditions of Theorem 8, for each  $i \geq 1$ ,

$$\sup_{t \leq T} \left| \lambda_n \int_0^t \beta_n(G_{n,i}(s)) ds - v_i(t) \right| \xrightarrow{P} 0,$$

where  $(g_i, v_i)$  is the unique solution of (7.3).

## 7.5 Properties of the Near Fixed Point

In this section we give some important properties of the near fixed point  $\mu_n$  that will be needed in the proofs of fluctuation theorems. Since  $\mu_n$  is defined in terms of the function  $\beta_n$ , we begin by giving some results on the asymptotic behavior of  $\beta_n$  and its derivatives. Proofs follow via elementary algebra and Taylor's approximation and can be found in Section 7.10. Roughly speaking, these results control the error between sampling with and without replacement of  $d_n$  servers from a collection of  $n$  servers. We first note that the function  $\beta_n$  is differentiable on  $(0, 1) \setminus \{\frac{d_n-1}{n}\}$  and the derivative is given as

$$\beta'_n(x) = \sum_{j=0}^{d_n-1} (1 - j/n)^{-1} \prod_{\substack{i=0 \\ i \neq j}}^{d_n-1} \frac{x - i/n}{1 - i/n} \text{ for } x \in (\frac{d_n-1}{n}, 1] \text{ and } \beta'_n(x) = 0 \text{ for } x \in (0, \frac{d_n-1}{n}). \quad (7.36)$$

As a convention, we set  $\beta'_n(x) = 0$  for  $x = \frac{d_n-1}{n}$ .

Note that  $f(t) = \frac{a+t}{b+t}$  is an increasing function of  $t$  on  $(-b, \infty)$  when  $b > a$ . Using this fact in (7.8) shows that, when  $d_n \leq n$ ,

$$0 \leq \beta_n(x) \leq x^{d_n} \doteq \gamma_n(x), \quad x \in [0, 1]. \quad (7.37)$$

Using the same fact in (7.36) shows that, for  $d_n < n$ ,

$$0 \leq \beta'_n(x) \leq \frac{d_n x^{d_n-1}}{1 - \frac{d_n}{n}}, x \in (0, 1). \quad (7.38)$$

The following lemma estimates the ratio between  $\beta_n$  and  $\gamma_n$  and its derivatives.

**Lemma 21.** *Assume  $\frac{d_n}{n} \rightarrow 0$ . Then for any  $\epsilon \in (0, 1)$ , as  $n \rightarrow \infty$ ,*

$$\sup_{x \in [\epsilon, 1]} \left| \frac{\beta'_n(x)/\beta_n(x)}{\gamma'_n(x)/\gamma_n(x)} - 1 \right| \rightarrow 0. \quad (7.39)$$

Furthermore, if  $\frac{d_n}{\sqrt{n}} \rightarrow 0$ , then

$$\sup_{x \in [\epsilon, 1]} \left| \frac{\beta_n(x)}{\gamma_n(x)} - 1 \right| \rightarrow 0 \quad \text{and} \quad \sup_{x \in [\epsilon, 1]} \left| \frac{\beta'_n(x)}{\gamma'_n(x)} - 1 \right| \rightarrow 0. \quad (7.40)$$

**Corollary 7.** *Assume  $d_n \ll n$ . Then for any  $\epsilon \in (0, 1)$*

$$\sup_{x \in [\epsilon, 1]} |\log \beta_n(x) - \log \gamma_n(x)| = O\left(\frac{d_n^2}{n}\right).$$

Recall the near fixed points  $\boldsymbol{\mu}_n = (\mu_{n,i})_{i \geq 1}$  introduced in Definition 9.

**Corollary 8.** *Suppose that  $d_n \ll n$ . Let  $i \in \mathbb{N}$  be such that  $\liminf_n \mu_{n,i} > 0$ . Then*

$$\lim_{n \rightarrow \infty} \frac{\lambda_n \mu_{n,i} \beta'_n(\mu_{n,i})}{d_n \mu_{n,i+1}} = 1.$$

**Lemma 22.** *Assume  $d_n \ll n$  and fix  $\epsilon \in (0, 1)$ . Then there is a  $C \in (0, \infty)$  and  $n_0 \in \mathbb{N}$  such that, if for some  $k \in \mathbb{N}$  and  $n_1 \in \mathbb{N}$ ,  $\mu_{n,k} \geq \epsilon$  for all  $n \geq n_1$ , then for all  $n \geq n_1 \vee n_0$*

$$\left| \log \mu_{n,k+1} - (\log \lambda_n) \left( \sum_{i=0}^k d_n^i \right) \right| \leq \frac{C}{n} \sum_{i=1}^k d_n^{i+1}.$$

**Corollary 9.** *Suppose that  $d_n \rightarrow \infty$  and that for some  $k \in \mathbb{N}$   $d_n^{k+1} \ll n$ . Suppose also that  $1 - \lambda_n = \frac{\xi_n + \log d_n}{d_n^k}$  where  $\xi_n \rightarrow -\log(\alpha) \in (-\infty, \infty]$  and  $\frac{\xi_n^2}{d_n} \rightarrow 0$ . Then  $\mu_{n,k} \rightarrow 1$  and  $\beta'_n(\mu_{n,k}) \rightarrow \alpha$ .*



**Lemma 23.** Suppose that  $\lambda_n \nearrow 1$ ,  $d_n \rightarrow \infty$  and  $d_n \ll n$ . Suppose also that, for some  $k \geq 2$ ,  $\mu_{n,k} \rightarrow 1$  and  $\beta'_n(\mu_{n,k}) \rightarrow \alpha \in [0, \infty)$  as  $n \rightarrow \infty$ . Then  $\beta'_n(\mu_{n,1}) \rightarrow \infty$  and for any  $i \in [k-1]$

$$\frac{\beta'_n(\mu_{n,i})}{\beta'_n(\mu_{n,1})} \rightarrow 1.$$

The following result is along the lines of Lemma 21. It allows for weaker assumptions on  $d_n$  but gives an approximation only in a neighborhood of 1.

**Lemma 24.** Suppose that  $\frac{d_n}{n^{2/3}} \rightarrow 0$ , as  $n \rightarrow \infty$ . Let  $\{\epsilon_n\}$  be a sequence in  $[0, 1]$  such that  $d_n \epsilon_n^2 \rightarrow 0$ . Then as  $n \rightarrow \infty$ :

$$\sup_{x \in [1-\epsilon_n, 1]} \left| \frac{\beta_n(x)}{\gamma_n(x)} - 1 \right| \rightarrow 0, \quad (7.41)$$

and

$$\sup_{x \in [1-\epsilon_n, 1]} \left| \frac{\beta'_n(x)}{\gamma'_n(x)} - 1 \right| \rightarrow 0. \quad (7.42)$$

The next result shows that if  $d_n \rightarrow \infty$ , then the behavior of  $\beta_n(x)$  is interesting only when  $x$  is sufficiently close to 1.

**Lemma 25.** Suppose that  $d_n \rightarrow \infty$ , and let  $\epsilon_n \doteq \frac{2 \log d_n}{d_n}$ . Then as  $n \rightarrow \infty$   $\sup_{x \in [0, 1-\epsilon_n]} |\beta_n(x)| \rightarrow 0$ . Furthermore, if  $\limsup_n \frac{d_n}{n} < 1$  then we also have  $\sup_{x \in [0, 1-\epsilon_n]} |\beta'_n(x)| \rightarrow 0$ .

## 7.6 Preliminary estimates under diffusion scaling

Recall the near fixed point  $\boldsymbol{\mu}_n$  from Definition 9 and the process  $\mathbf{Z}_n$  introduced in (7.1). Also, recall the maps  $\mathbf{a}_n$  and  $\mathbf{b}$  from Remark 7.3.1. We will extend the definition of  $\beta_n$  and  $\beta'_n$  to  $\mathbb{R}$  by setting  $\beta_n(x) = \beta'_n(x) = 0$  for  $x < 0$ . Further, in what follows, for  $z < 0$  and real valued integrable function  $h(\cdot)$ , the integral  $\int_{[0, z]} h(u) du = - \int_{[z, 0]} h(u) du$ . We start by giving a semimartingale decomposition for  $\mathbf{Z}_n$ .

**Lemma 26.** For  $t \geq 0$ ,  $\mathbf{Z}_n(t)$  satisfies

$$\mathbf{Z}_n(t) = \mathbf{Z}_n(0) + \int_0^t \mathbf{A}_n(\mathbf{Z}_n(s)) ds - \int_0^t \mathbf{b}(\mathbf{Z}_n(s)) ds + \sqrt{n} \mathbf{M}_n(t) \quad (7.43)$$

where  $\mathbf{A}_n : l_\infty \rightarrow l_\infty$  is given as

$$\mathbf{A}_n(\mathbf{z})_i = t_{n,i-1}(z_{i-1}) - t_{n,i}(z_i), \quad i \in \mathbb{N} \quad (7.44)$$

and for  $i \in \mathbb{N}$

$$\begin{aligned} t_{n,i}(z) &\doteq \lambda_n \int_{[0,z]} \beta'_n(\mu_{n,i} + y/\sqrt{n}) dy, \quad z \in \mathbb{R} \\ t_{n,0}(z) &\doteq 0, \quad z \in \mathbb{R}. \end{aligned} \quad (7.45)$$

*Proof.* From (7.23) and since  $\mathbf{a}_n(\boldsymbol{\mu}_n) = \mathbf{b}(\boldsymbol{\mu}_n)$ ,

$$\begin{aligned} \sqrt{n}(\mathbf{G}_n(t) - \boldsymbol{\mu}_n) &= \sqrt{n}(\mathbf{G}_n(0) - \boldsymbol{\mu}_n) + \int_0^t \sqrt{n}\{\mathbf{a}_n(\mathbf{G}_n(s)) - \mathbf{a}_n(\boldsymbol{\mu}_n)\} ds \\ &\quad - \int_0^t \sqrt{n}\{\mathbf{b}(\mathbf{G}_n(s)) - \mathbf{b}(\boldsymbol{\mu}_n)\} ds + \sqrt{n}\mathbf{M}_n(t) \end{aligned}$$

Let  $\mathbf{A}_n(\mathbf{z}) \doteq \sqrt{n}\{\mathbf{a}_n(\boldsymbol{\mu}_n + n^{-1/2}\mathbf{z}) - \mathbf{a}_n(\boldsymbol{\mu}_n)\}$ . By the definition of  $\mathbf{a}_n$  we see that (7.44) holds where

$$t_{n,i}(z) \doteq \begin{cases} \lambda_n \sqrt{n}\{\beta_n(\mu_{n,i} + z/\sqrt{n}) - \beta_n(\mu_{n,i})\} & \text{for } i \geq 1 \\ 0 & \text{if } i = 0 \end{cases} \quad (7.46)$$

Clearly, the  $t_{n,i}$  defined in (7.46) is same as that given in (7.45). The result follows.  $\square$

**Lemma 27.** Suppose that  $d_n \rightarrow \infty$ ,  $\lambda_n \rightarrow 1$ , and for some  $k \geq 1$ ,  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{f}_k$  in  $l_1$ . Then there is a standard Brownian motion  $B$  so that  $\sqrt{n}\mathbf{M}_n \Rightarrow \sqrt{2}B\mathbf{e}_k$  in  $\mathbb{D}([0, \infty) : l_2)$ .

*Proof.* Fix  $T > 0$ . Since  $\mathbf{G}_n(0) \rightarrow \mathbf{f}_k$  and  $\mathbf{f}_k$  is a fixed point of (7.3), by Theorem 8,  $\mathbf{G}_n \xrightarrow{P} \mathbf{f}_k$  in  $\mathbb{D}([0, T] : l_1)$ , where  $\mathbf{f}_k$  here is viewed as the function on  $[0, T]$  that takes the constant value  $\mathbf{f}_k \in l_1$ . Moreover, by Remark 7.4.1, for every  $i \geq 1$   $V_{n,i}(t) \doteq \lambda_n \int_0^t \beta_n(G_{n,i}(s)) ds$  converge uniformly on  $[0, T]$  in probability to  $v_i(t)$ , where  $v_i$  solves

$$v_i = \hat{\Gamma}_1(f_{k,i} - (f_{k,i} - f_{k,i+1})\text{id} + v_{i-1}(\cdot)), \quad i \geq 1, \quad (7.47)$$

and  $v_0(t) \doteq t$ , where recall that  $\text{id} : [0, T] \rightarrow [0, T]$  is the identity map. Recalling the definition of  $f_k$  we see by a recursive argument that

$$v_i(t) \doteq \begin{cases} t & \text{if } i < k \\ 0 & \text{if } i \geq k. \end{cases} \quad (7.48)$$

Combining this with (7.22), we have for each  $i \geq 1$

$$\begin{aligned} \langle \sqrt{n}M_{n,i} \rangle_\cdot &= \int_0^\cdot (G_{n,i}(s) - G_{n,i+1}(s))ds + \lambda_n \int_0^\cdot (\beta_n(G_{n,i-1}(s)) - \beta_n(G_{n,i}(s))) ds \\ &\rightarrow (f_{k,i} - f_{k,i+1})\text{id} + v_{i-1}(\cdot) - v_i(\cdot) = H(\cdot), \end{aligned}$$

in probability in  $C([0, T] : \mathbb{R})$  where

$$H(t) \doteq \begin{cases} 2t & \text{if } i = k \\ 0 & \text{if } i \neq k \end{cases}, \quad t \in [0, T].$$

Adding (7.22) over  $i$ , we have for  $t \in [0, T]$ ,

$$\sum_{i>k} \langle \sqrt{n}M_{n,i} \rangle_t \leq \int_0^t G_{n,k+1}(s)ds + \lambda_n \int_0^t \beta_n(G_{n,k})(s)ds. \quad (7.49)$$

The process on the right side converges in probability in  $C([0, T] : \mathbb{R})$  to  $f_{k,k+1}\text{id} + v_k(\cdot) = 0$  and thus  $\sum_{i>k} \langle \sqrt{n}M_{n,i} \rangle_T$  converges to 0 in probability. By Doob's maximal inequality,

$$n\mathbf{E} \sup_{t \leq T} \sum_{i>k} M_{n,i}^2(t) \leq 4\mathbf{E} \sum_{i>k} \langle \sqrt{n}M_{n,i} \rangle_T \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

where the last convergence follows by the dominated convergence theorem on noting that the right side of (7.49) is bounded above by  $\sup_n(1 + \lambda_n) < \infty$ . The result now follows on using the martingale central limit theorem (cf. (46, Theorem 7.1.4)) for the  $k$ -dimensional martingale sequence  $(\sqrt{n}M_{n,1}, \dots, \sqrt{n}M_{n,k})$ .  $\square$

Recall the functions  $t_{n,i}$  from Lemma 26.

**Lemma 28.** Assume that for some  $r \in \mathbb{N}$ ,  $\limsup_{n \rightarrow \infty} \mu_{n,r} < 1$ . Then for any  $L > 0$

$$\limsup_{n \rightarrow \infty} \sup_{i \geq r} \sup_{0 < |z| \leq L} \left| \frac{t_{n,i}(z)}{z} \right| = 0$$

*Proof.* By (7.45):

$$\begin{aligned} \sup_{i \geq r} \sup_{0 < |z| \leq L} \left| \frac{t_{n,i}(z)}{z} \right| &\leq \lambda_n \sup_{i \geq r} \sup_{0 < |z| \leq L} \sup_{|y| \leq z} \left| \beta'_n \left( \mu_{n,i} + \frac{y}{\sqrt{n}} \right) \right| \\ &= \lambda_n \sup_{i \geq r} \sup_{|z| \leq L} \left| \beta'_n \left( \mu_{n,i} + \frac{z}{\sqrt{n}} \right) \right| \leq \lambda_n \sup_{0 \leq x \leq \mu_{n,r} + \frac{L}{\sqrt{n}}} \beta'_n(x) \end{aligned}$$

which converges to 0 by Lemma 25, since  $\limsup_{n \rightarrow \infty} \left( \mu_{n,r} + \frac{L}{\sqrt{n}} \right) < 1$ .  $\square$

For  $L \in (0, \infty)$  define the stopping time

$$\tau_{n,L} \doteq \inf \left\{ t \mid \|\mathbf{Z}_n(t)\|_2 \geq L - \frac{1}{\sqrt{n}} \right\}. \quad (7.50)$$

Since the jumps of  $\mathbf{Z}_n$  are of size  $\frac{1}{\sqrt{n}}$ , we see that, for any  $T > 0$

$$\|\mathbf{Z}_n\|_{2,T \wedge \tau_{n,L}} \leq L. \quad (7.51)$$

Recall from Section 7.1.2 the vector  $\mathbf{z}_{r+} \in \mathbb{R}^\infty$  associated with a vector  $\mathbf{z} \in \mathbb{R}^\infty$ .

**Lemma 29.** Suppose that as  $n \rightarrow \infty$ ,  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{f}_k$  in  $l_1^\downarrow$  and  $\mathbf{Z}_{n,r+}(0) \xrightarrow{P} \mathbf{0}$  in  $l_2$  for some  $r > k$ . Then for any  $T, L > 0$ ,  $\|\mathbf{Z}_{n,r+}\|_{2,T \wedge \tau_{n,L}} \xrightarrow{P} 0$ .

*Proof.* For  $i > k$  and  $z \in \mathbb{R}$ , let  $\Delta_{n,i}(z) \doteq \frac{t_{n,i}(z)}{z} \mathbb{I}_{\{z \neq 0\}}$ . Then, since  $\lim_{n \rightarrow \infty} \mu_{n,k+1} = 0$ , by Lemma

28

$$\delta_{n,L} \doteq \sup_{i \geq k+1} \sup_{|z| \leq L} |\Delta_{n,i}(z)| \rightarrow 0, \text{ as } n \rightarrow \infty. \quad (7.52)$$

Next, from (7.43), for  $i \geq r+1 > k+1$

$$\begin{aligned} Z_{n,i}(t \wedge \tau_n) &= Z_{n,i}(0) + \int_0^{t \wedge \tau_n} \Delta_{n,i-1}(Z_{n,i-1}(s)) Z_{n,i-1}(s) ds - \int_0^{t \wedge \tau_n} \Delta_{n,i}(Z_{n,i}(s)) Z_{n,i}(s) ds \\ &\quad - \int_0^{t \wedge \tau_n} (Z_{n,i}(s) - Z_{n,i+1}(s)) ds + \sqrt{n} M_{n,i}(t \wedge \tau_n) \end{aligned}$$

where we use  $\tau_n$  instead of  $\tau_{n,L}$  for notational simplicity. Then, observing from (7.52) that  $\sup_{i \geq k+1} \sup_{t \in [0, \tau_n]} |\Delta_{n,i}(Z_{n,i}(t))| \leq \delta_{n,L}$ , we have

$$\begin{aligned} |Z_{n,i}(t \wedge \tau_n)| &\leq |Z_{n,i}(0)| + \delta_{n,L} \int_0^{t \wedge \tau_n} (|Z_{n,i-1}(s)| + |Z_{n,i}(s)|) ds \\ &\quad + \int_0^{t \wedge \tau_n} (|Z_{n,i}(s)| + |Z_{n,i+1}(s)|) ds + |\sqrt{n} M_{n,i}(t \wedge \tau_n)|. \end{aligned} \quad (7.53)$$

Define maps  $\mathbf{A}_1, \mathbf{A}_2 : \mathbb{R}^\infty \rightarrow \mathbb{R}^\infty$  by

$$\begin{aligned} (\mathbf{A}_1 \mathbf{x})_i &= \begin{cases} x_1 & i = 1 \\ x_{i-1} + x_i & i \geq 2 \end{cases} \\ (\mathbf{A}_2 \mathbf{x})_i &= x_i + x_{i+1}, \quad i \in \mathbb{N}. \end{aligned}$$

Then by collecting (7.53) over all  $i \geq r+1$  we get

$$\begin{aligned} |\mathbf{Z}_{n,r+}(t \wedge \tau_n)| &\leq |\mathbf{Z}_{n,r+}(0)| + \delta_{n,M} \int_0^{t \wedge \tau_n} \mathbf{A}_1 |\mathbf{Z}_{n,r+}(s)| ds + \delta_{n,M} \int_0^{t \wedge \tau_n} |Z_{n,r}(s)| \mathbf{e}_1 ds \\ &\quad + \int_0^{t \wedge \tau_n} \mathbf{A}_2 |\mathbf{Z}_{n,r+}(s)| ds + |\sqrt{n} \mathbf{M}_{n,r+}(t \wedge \tau_n)| \end{aligned} \quad (7.54)$$

where the absolute values and the integrals are interpreted as being coordinate-wise for infinite dimensional vectors. Now noting that the maps  $\mathbf{A}_i$ , when considered from  $l_2 \rightarrow l_2$ , are bounded linear operators with norm bounded by 2, we have for  $i = 1, 2$ ,

$$\left\| \int_0^{t \wedge \tau_n} \mathbf{A}_i |\mathbf{Z}_{n,r+}(s)| ds \right\|_2 \leq \int_0^{t \wedge \tau_n} 2 \|\mathbf{Z}_{n,r+}(s)\|_2 ds.$$

Using the triangle inequality in (7.54) shows for any  $t \leq T$

$$\begin{aligned} \|\mathbf{Z}_{n,r+}(t \wedge \tau_n)\|_2 &\leq \|\mathbf{Z}_{n,r+}(0)\|_2 + \|\sqrt{n}\mathbf{M}_{n,r+}\|_{2,T} + \delta_{n,M}MT \\ &\quad + 2(1 + \delta_{n,M}) \int_0^{t \wedge \tau_n} \|\mathbf{Z}_{n,r+}(s)\|_2 ds \end{aligned}$$

where we have used that  $\int_0^{t \wedge \tau_n} |Z_{n,r}(s)| ds \leq Lt$ . Hence, using Gronwall's inequality

$$\|\mathbf{Z}_{n,r+}\|_{2,T \wedge \tau_n} \leq \left( \|\mathbf{Z}_{n,r+}(0)\|_2 + \delta_{n,L}LT + \|\sqrt{n}\mathbf{M}_{n,r+}\|_{2,T} \right) e^{2(1+\delta_{n,L})T}$$

Now, as  $n \rightarrow \infty$ ,  $\|\mathbf{Z}_{n,r+}(0)\|_2 \xrightarrow{P} 0$  by assumption,  $\delta_{n,L} \rightarrow 0$  by (7.52), and  $\|\sqrt{n}\mathbf{M}_{n,r+}\|_{2,T} \xrightarrow{P} 0$  by Lemma 27. The result follows.  $\square$

The following elementary lemma will allow us to replace  $\tau_{n,L} \wedge T$  with  $T$  in various convergence results. The proof is omitted.

**Lemma 30.** *Fix  $T \in [0, \infty)$ . Suppose for each  $n \in \mathbb{N}$  and  $L > 0$  that  $\tau_{n,L}$  is a  $[0, T]$  valued random variable such that  $\lim_{L \rightarrow \infty} \sup_n \mathbf{P}(\tau_{n,L} < T) \rightarrow 0$  for some  $T > 0$ . Suppose that there is a sequence of stochastic processes  $\{F_n\}_{n \in \mathbb{N}}$  with sample paths in  $\mathbb{D}([0, T] : \mathbb{R})$  such that for each  $L > 0$   $|F_n|_{*, T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . Then in fact  $|F_n|_{*, T} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ .*

The next lemma gives conditions under which the near fixed point  $\boldsymbol{\mu}_n$  converges to  $\mathbf{f}_1$ .

**Lemma 31.** *Let  $0 \leq \epsilon_n \doteq 1 - \lambda_n$  be such that  $\epsilon_n \rightarrow 0$  and  $\epsilon_n d_n \rightarrow \infty$ . Then  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_1$  in  $l_1$  as  $n \rightarrow \infty$ .*

*Proof.* Since  $d_n \rightarrow \infty$  under our assumptions, in order to show  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_1$  in  $l_1$  it suffices to show that (1)  $\mu_{n,1} \rightarrow 1$ , and (2)  $\mu_{n,2} \rightarrow 0$ . The convergence in (1) is immediate on observing that  $\mu_{n,1} = \lambda_n = 1 - \epsilon_n \rightarrow 1$ , and (2) follows by noting from Definition 9 and (7.37) that  $\mu_{n,2} \leq \mu_{n,1}^{d_n} = (1 - \epsilon_n)^{d_n} \leq e^{-\epsilon_n d_n} \rightarrow 0$ .  $\square$

The following lemma gives a convenient approximation of the term  $t_{n,1}$  introduced in (7.45) in terms of certain exponentials.

**Lemma 32.** Suppose  $d_n \rightarrow \infty$  and  $d_n \ll n^{2/3}$ . Let  $\lambda_n = 1 - \left( \frac{\log d_n}{d_n} + \frac{\alpha_n}{\sqrt{n}} \right)$  for some real sequence  $\{\alpha_n\}$  satisfying  $\frac{d_n \alpha_n^2}{n} \rightarrow 0$ . Then, for any  $L > 0$ ,

$$\limsup_{n \rightarrow \infty} \sup_{0 < |z| \leq L} \left| \frac{\exp\left(\frac{d_n}{\sqrt{n}}(z - \alpha_n)\right) - \exp\left(-\frac{d_n}{\sqrt{n}}\alpha_n\right)}{t_{n,1}(z)d_n/\sqrt{n}} - 1 \right| = 0. \quad (7.55)$$

*Proof.* We only consider the case  $0 < z \leq L$ . The case  $-L \leq z < 0$  is treated similarly. Recall that  $\mu_{n,1} = \lambda_n$ . Noting that  $d_n(1 - \lambda_n + \frac{L}{\sqrt{n}})^2 \leq 4d_n(\frac{\log^2 d_n}{d_n^2} + \alpha_n^2/n + L^2/n) \rightarrow 0$  we have on applying Lemma 24 with  $\epsilon_n = (1 - \lambda_n + \frac{L}{\sqrt{n}})$  that, for any  $|z| \leq L$ ,

$$\begin{aligned} t_{n,1}(z) &= (1 + o(1)) \int_0^z \gamma'_n\left(\lambda_n + \frac{y}{\sqrt{n}}\right) dy \\ &= (1 + o(1)) \int_0^z \exp\left((d_n - 1) \log\left\{\lambda_n + \frac{y}{\sqrt{n}}\right\} + \log d_n\right) dy \\ &= (1 + o(1)) \int_0^z \exp\left(d_n \log\left\{\lambda_n + \frac{y}{\sqrt{n}}\right\} + \log d_n\right) dy \end{aligned}$$

Using expansion for  $\log(1 + h)$  around  $h = 0$  and once more the fact that  $d_n\left(1 - \lambda_n + \frac{L}{\sqrt{n}}\right)^2 \rightarrow 0$ ,

$$\begin{aligned} t_{n,1}(z) &= (1 + o(1)) \int_0^z \exp\left(d_n \left\{\lambda_n - 1 + \frac{y}{\sqrt{n}}\right\} + \log d_n\right) dy \\ &= (1 + o(1)) \int_0^z \exp\left(\frac{d_n}{\sqrt{n}}(y - \alpha_n)\right) dy \\ &= (1 + o(1)) \frac{\exp\left(\frac{d_n}{\sqrt{n}}(z - \alpha_n)\right) - \exp\left(-\frac{d_n}{\sqrt{n}}\alpha_n\right)}{d_n/\sqrt{n}} \end{aligned}$$

which proves (7.55). □

Proof of the following lemma proceeds by standard arguments but we provide details in Section 7.10.9.

**Lemma 33.** Fix  $T > 0$ . Let  $f, g, M$  be three bounded measurable functions from  $[0, T] \rightarrow \mathbb{R}$  and assume further that  $M$  is a right continuous bounded variation function. Suppose that  $m \doteq \inf_{s \in [0, T \wedge \tau]} f(s) > 0$  for some  $\tau \geq 0$ . Let  $z : [0, T] \rightarrow \mathbb{R}$  be a bounded measurable function that

satisfies for every  $t \in [0, T]$

$$z(t) = z(0) - \int_0^t f(s)z(s)ds + \int_0^t g(s)ds + M(t). \quad (7.56)$$

Then for any  $t \in [0, T \wedge \tau]$

$$|z(t)| \leq \frac{|g|_{*, T \wedge \tau}}{m} + 2|M|_{*, T \wedge \tau} + e^{-mt}(|z(0)| + |M(0)|).$$

**Lemma 34.** Fix  $T \in (0, \infty)$ . For each  $n$ , let  $V_n$  be a martingale with respect to some filtration  $\{\mathcal{G}_t^n\}$  such that  $V_n(0) = 0$ . Let  $(r_n)_{n=1}^\infty$  be a positive sequence so that  $\lim_{n \rightarrow \infty} r_n = +\infty$ . Suppose that there is a  $C \in (0, \infty)$  such that for all  $n \in \mathbb{N}$  and  $t \in [0, T]$ ,  $\langle V_n \rangle_t \leq Ct$ . Then for any  $\epsilon > 0$

$$\mathbf{P}\left(\sup_{t \leq T} (V_n(t) - r_n t) > \epsilon\right) \rightarrow 0$$

as  $n \rightarrow \infty$ .

*Proof.* Let  $\delta_n \doteq \frac{1}{\sqrt{r_n}}$ . Then

$$\begin{aligned} \mathbf{P}\left(\sup_{0 \leq t \leq T} [V_n(t) - r_n t] > \epsilon\right) &\leq \mathbf{P}\left(\sup_{0 \leq t \leq \delta_n} |V_n(t)| > \epsilon\right) + \mathbf{P}\left(\sup_{\delta_n < t \leq T} |V_n(t)| > r_n \delta_n\right) \\ &\leq \frac{4\mathbf{E}V_n(\delta_n)^2}{\epsilon^2} + \frac{4\mathbf{E}V_n(T)^2}{(r_n \delta_n)^2} \\ &= \frac{4\mathbf{E}\langle V_n \rangle_{\delta_n}}{\epsilon^2} + \frac{4\mathbf{E}\langle V_n \rangle_T}{(r_n \delta_n)^2} \leq \frac{4C\delta_n}{\epsilon^2} + \frac{4CT}{(r_n \delta_n)^2} \rightarrow 0 \end{aligned}$$

where the inequality on the second line is from Doob's maximal inequality.  $\square$

## 7.7 Proof of Theorem 9

Now we start with some preliminary lemmas. Recall from Remark 7.2.3(ii) that under the hypothesis of Theorem 9 we have  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_k \in l_1^\downarrow$  as  $n \rightarrow \infty$ . Along with the tightness of  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$  this also shows that  $\mathbf{G}_n(0) \rightarrow \mathbf{f}_k \in l_1^\downarrow$  as  $n \rightarrow \infty$ .



**Lemma 35.** Let  $d_n \rightarrow \infty$ ,  $\frac{d_n}{\sqrt{n}} \rightarrow 0$ , and  $\lambda_n \nearrow 1$ . Assume that for some  $k \in \mathbb{N}$ ,  $\mu_n \rightarrow \mathbf{f}_k$  in  $\mathbb{R}^\infty$ . Then for any  $M > 0$  and  $1 \leq i \leq k$ , as  $n \rightarrow \infty$

$$\sup_{0 < |z| \leq M} \left| (\beta'_n(\mu_{n,i})z)^{-1} \lambda_n \int_0^z \beta'_n(\mu_{n,i} + y/\sqrt{n}) dy - 1 \right| \rightarrow 0 \quad (7.57)$$

*Proof.* To prove (7.57), we will approximate  $\beta'_n(x)$  by  $\gamma'_n(x)$ . Using Lemma 21

$$\epsilon_n \doteq \sup_{x \in [1/2, 1]} \left| \frac{\beta'_n(x)}{\gamma'_n(x)} - 1 \right| \rightarrow 0.$$

Since  $\mu_{n,k} \rightarrow 1$ , there is an  $N_0$  so that for  $n \geq N_0$ ,  $\mu_{n,i} + \frac{y}{\sqrt{n}} \geq \frac{1}{2}$ , for any  $i \leq k$  and  $y \in \mathbb{R}$  with  $|y| \leq L$ . Hence uniformly in  $0 < |z| \leq L$  and  $i \leq k$ :

$$\begin{aligned} \frac{\lambda_n}{z} \int_0^z \frac{\beta'_n(\mu_{n,i} + \frac{y}{\sqrt{n}})}{\beta'_n(\mu_{n,i})} dy &= \frac{1 + o(1)}{z} \int_0^z \frac{\gamma'_n(\mu_{n,i} + \frac{y}{\sqrt{n}})}{\gamma'_n(\mu_{n,i})} dy \\ &= \frac{1 + o(1)}{z} \int_0^z \left( 1 + \frac{y}{\sqrt{n}\mu_{n,i}} \right)^{d_n-1} dy \\ &= \frac{1 + o(1)}{z} \int_0^z \exp \left\{ (d_n - 1) \log \left( 1 + \frac{y}{\sqrt{n}\mu_{n,i}} \right) \right\} dy \\ &= \frac{1 + o(1)}{z} \int_0^z \exp \left\{ O \left( \frac{d_n L}{\sqrt{n}\mu_{n,k}} \right) \right\} dy \rightarrow 1 \end{aligned}$$

This shows (7.57). □

**Remark 7.7.1.** Suppose that the hypothesis of Lemma 35 hold. Recall the definition of  $\Delta_{n,i}$  for  $i > k$  from the proof of Lemma 29. We extend this definition by setting

$$\Delta_{n,i}(z) \doteq t_{n,i}(z)/(\beta'_n(\mu_{n,i})z) \mathbb{I}_{\{z \neq 0\}} - 1 \text{ if } 1 \leq i \leq k \quad (7.58)$$

where  $t_{n,i}$  is defined by (7.45). With this extension

$$t_{n,i}(z) = \begin{cases} \beta'_n(\mu_{n,i})(1 + \Delta_{n,i}(z))z & \text{if } 1 \leq i \leq k \\ \Delta_{n,i}(z)z & \text{if } i > k \end{cases}. \quad (7.59)$$

Using this notation, Lemma 35 and Lemma 28 show that, for any  $L > 0$

$$\gamma_{n,L} \doteq \sup_{i \in \mathbb{N}} \sup_{0 < |z| \leq L} |\Delta_{n,i}(z)| \rightarrow 0 \text{ as } n \rightarrow \infty. \quad (7.60)$$

The following corollary is an immediate consequence of Remark 7.7.1 and Lemma 26.

**Corollary 10.** *Under the hypothesis of Lemma 35,  $\mathbf{Z}_n$  satisfies the following integral equations.*

For  $i = 1$

$$Z_{n,1}(t) = Z_{n,1}(0) - \int_0^t \beta'_n(\mu_{n,1})(1 + \Delta_{n,1}(Z_{n,1}(s)))Z_{n,1}(s)ds - \int_0^t (Z_{n,1}(s) - Z_{n,2}(s))ds + \sqrt{n}M_{n,1}(t)$$

For  $i \in \{2, \dots, k\}$

$$\begin{aligned} Z_{n,i}(t) = & Z_{n,i}(0) + \int_0^t \beta'_n(\mu_{n,i-1})(1 + \Delta_{n,i-1}(Z_{n,i-1}(s)))Z_{n,i-1}(s)ds \\ & - \int_0^t \beta'_n(\mu_{n,i})(1 + \Delta_{n,i}(Z_{n,i}(s)))Z_{n,i}(s)ds - \int_0^t (Z_{n,i}(s) - Z_{n,i+1}(s))ds + \sqrt{n}M_{n,i}(t). \end{aligned}$$

For  $i = k + 1$

$$\begin{aligned} Z_{n,k+1}(t) = & Z_{n,k+1}(0) + \int_0^t \beta'_n(\mu_{n,k})(1 + \Delta_{n,k}(Z_{n,k}(s)))Z_{n,k}(s)ds \\ & - \int_0^t \Delta_{n,k+1}(Z_{n,k+1}(s))Z_{n,k+1}(s) - \int_0^t (Z_{n,k+1}(s) - Z_{n,k+2}(s))ds + \sqrt{n}M_{n,k+1}(t), \end{aligned}$$

For  $i > k + 1$

$$\begin{aligned} Z_{n,i}(t) = & Z_{n,i}(0) + \int_0^t \Delta_{n,i-1}(Z_{n,i-1}(s))Z_{n,i-1}(s)ds - \int_0^t \Delta_{n,i}(Z_{n,i}(s))Z_{n,i}(s)ds \\ & - \int_0^t (Z_{n,i}(s) - Z_{n,i+1}(s))ds + \sqrt{n}M_{n,i}(t), \end{aligned}$$

where  $\Delta_{n,i}$  is as in Remark 7.7.1.

Finally, if  $Y_{n,1} \doteq \sum_{i=1}^k Z_{n,i}$ , then

$$\begin{aligned} Y_{n,1}(t) &= Y_{n,1}(0) - \int_0^t \beta'_n(\mu_{n,k})(1 + \Delta_{n,k}(Z_{n,k}(s)))Z_{n,k}(s)ds \\ &\quad - \int_0^t (Z_{n,1}(s) - Z_{n,k+1}(s))ds + \sum_{i=1}^k \sqrt{n}M_{n,i}(t) \end{aligned} \quad (7.61)$$

**Lemma 36.** Suppose  $\lambda_n \nearrow 1$ ,  $d_n \rightarrow \infty$  and  $d_n \ll n$ . Assume that for some  $k \geq 2$   $\mu_{n,k} \rightarrow 1$  and  $\beta'_n(\mu_{n,k}) \rightarrow \alpha \in [0, \infty)$  as  $n \rightarrow \infty$ . Define the  $k-1 \times k-1$  tridiagonal matrix  $A_n(s)$  as

$$\begin{aligned} A_n(s)[j, j] &= a_{n,j}(s) + 1, \quad 1 \leq j \leq k-1, \\ A_n(s)[j, j+1] &= -1, \quad 1 \leq j \leq k-2, \\ A_n(s)[j, j-1] &= -a_{n,j-1}(s), \quad 2 \leq j \leq k-1, \end{aligned} \quad (7.62)$$

and for all other  $j, k$ ,  $A_n(s)[j, k] = 0$ , where  $a_{n,i}(s) \doteq \beta'_n(\mu_{n,i})(1 + \Delta_{n,i}(Z_{n,i}(s)))$ . Then for any  $T, L \in (0, \infty)$

$$\lim_{n \rightarrow \infty} \inf_{s \in [0, T \wedge \tau_{n,L}]} \inf_{\vec{x} \in \mathbb{R}^{k-1} \setminus \{0\}} \frac{\vec{x}^t A_n(s) \vec{x}}{\|\vec{x}\|^2} = +\infty.$$

*Proof.* Let  $b_{n,i}(s) \doteq a_{n,i}(s) + 1$ . and  $B_n(s) \doteq A_n(s) + A_n(s)^t$ . Then  $B_n(s)$  is a symmetric tridiagonal matrix with entries

$$\begin{aligned} B_n(s)[j, j] &= 2b_{n,j}(s), \quad 1 \leq j \leq k-1, \\ B_n(s)[j, j+1] &= -b_{n,j}, \quad 1 \leq j \leq k-2, \\ B_n(s)[j, j-1] &= -b_{n,j-1}(s), \quad 2 \leq j \leq k-1, \end{aligned} \quad (7.63)$$

Let  $b_n \doteq \beta'_n(\mu_{n,1})$ . By Lemma 23,  $b_n \rightarrow \infty$  and by the uniform convergence in (7.60) and Lemma 23 once more

$$\max_{i \leq k-1} \sup_{s \in [0, T \wedge \tau_{n,L}]} \left| \frac{b_{n,i}(s)}{b_n} - 1 \right| \rightarrow 0 \text{ as } n \rightarrow \infty.$$

This in particular shows that

$$\sup_{s \in [0, T \wedge \tau_{n,L}]} \left\| \frac{1}{b_n} B_n(s) - H \right\|_F \rightarrow 0, \quad (7.64)$$

where  $\|\cdot\|_F$  is the Frobenius norm and  $H$  is the  $k-1 \times k-1$  tridiagonal matrix given as

$$\begin{aligned} H[j, j] &= 2, \quad 1 \leq j \leq k-1, \\ H[j, j+1] &= -1, \quad 1 \leq j \leq k-2, \\ H[j, j-1] &= -1, \quad 2 \leq j \leq k-1, \end{aligned}$$

Note for any  $\vec{x} = (x_1, x_2, \dots, x_{k-1}) \in \mathbb{R}^{k-1}$  by completing squares

$$\vec{x}^t H \vec{x} = x_1^2 + (x_2 - x_1)^2 + (x_3 - x_2)^2 + \dots + (x_{k-2} - x_{k-1})^2 + x_{k-1}^2,$$

which is always non-zero if  $\vec{x} \neq 0$ . Let  $c \doteq \inf_{\|\vec{x}\|=1} \vec{x}^t H \vec{x}$ . Since the unit sphere is compact, the minimum is attained and hence  $c > 0$ . This shows that  $H$  is a positive definite matrix.

Finally,

$$\begin{aligned} \vec{x}^t \frac{1}{b_n} B_n(s) \vec{x} &= \vec{x}^t H \vec{x} + \vec{x}^t \left( \frac{1}{b_n} B_n - H \right) \vec{x} \\ &\geq \vec{x}^t H \vec{x} - \|b_n^{-1} B_n - H\|_F \|\vec{x}\|^2 \geq (c - \|b_n^{-1} B_n - H\|_F) \|\vec{x}\|^2. \end{aligned}$$

By (7.64), there is an  $N_0 \in \mathbb{N}$  so that for each  $n \geq N_0$  and each  $s \in [0, T \wedge \tau_{n,L}]$ ,  $\|b_n^{-1} B_n - H\|_F \leq c/2$ . Hence for each  $\vec{x} \in \mathbb{R}^{k-1}$

$$2\vec{x}^t A_n(s) \vec{x} = \vec{x}^2 B_n(s) \vec{x} \geq (c/2) b_n \|\vec{x}\|^2.$$

Since  $b_n \rightarrow \infty$ , this completes the proof.  $\square$

**Lemma 37.** *Suppose that the hypothesis of Theorem 9 holds with  $k \geq 2$  and let  $\vec{X}_n \doteq (Z_{n,1}, Z_{n,2}, \dots, Z_{n,k-1})$ . Then for  $L, T, \epsilon \in (0, \infty)$*

$$P \left( \sup_{s \in [0, T \wedge \tau_{n,L}]} \|\vec{X}_n(s)\| > \|\vec{X}_n(0)\| + \epsilon \right) \rightarrow 0, \quad (7.65)$$

and

$$\sup_{s \in [\epsilon, T \wedge \tau_{n,L}]} \|\vec{X}_n(s)\| \xrightarrow{P} 0, \quad (7.66)$$

as  $n \rightarrow \infty$ .

*Proof.* Let  $\vec{W}_n = (\sqrt{n}M_{n,1}, \dots, \sqrt{n}M_{n,k-1})$ . Then by Corollary 10

$$\vec{X}_n(t) = \vec{X}_n(0) - \int_0^t A_n(s) \vec{X}_n(s) ds + \vec{e}_{k-1} \int_0^t Z_{n,k}(s) ds + \vec{W}_n(t). \quad (7.67)$$

where  $\vec{e}_{k-1}$  is the vector  $(0, 0, \dots, 0, 1)' \in \mathbb{R}^{k-1}$  and  $A_n(s)$  is  $k-1 \times k-1$  matrix defined in (7.62). Using Ito's formula to the function  $f(\vec{x}) = \|\vec{x}\|^2$  along with the semimartingale representation from (7.67)

$$\begin{aligned} \|\vec{X}_n(t)\|^2 &= \|\vec{X}_n(0)\|^2 + 2 \int_{0+}^t \langle \vec{X}_n(s-), d\vec{X}_n(s) \rangle + [\vec{W}_n]_t \\ &= \|\vec{X}_n(0)\|^2 - 2 \int_0^t \langle \vec{X}_n(s), A_n(s) \vec{X}_n(s) \rangle ds + 2 \int_0^t Z_{n,k}(s) \langle \vec{X}_n(s), \vec{e}_{k-1} \rangle ds \\ &\quad + 2 \int_0^t \langle \vec{X}_n(s-), d\vec{W}_n(s) \rangle + [\vec{W}_n]_t, \end{aligned} \quad (7.68)$$

where  $[\vec{W}_n]_t \doteq \sum_{i=1}^{k-1} [\sqrt{n}M_{n,i}]_t$ . Define

$$\begin{aligned} f_n(s) &\doteq \frac{\vec{X}_n(s)^t A_n(s) \vec{X}_n(s)}{\|\vec{X}_n(s)\|^2} \mathbb{I}_{\{\vec{X}_n(s) \neq 0\}} + n \mathbb{I}_{\{\vec{X}_n(s) = 0\}} \\ g_n(s) &\doteq 2Z_{n,k}(s)Z_{n,k-1}(s) \\ B_n(s) &\doteq 2 \int_0^s \langle \vec{X}_n(s-), d\vec{W}_n(s) \rangle + [\vec{W}_n]_s \end{aligned}$$

then (7.68) becomes

$$\|\vec{X}_n(t)\|^2 = \|\vec{X}_n(0)\|^2 - 2 \int_0^t f_n(s) \|\vec{X}_n(s)\|^2 ds + \int_0^t g_n(s) ds + B_n(t). \quad (7.69)$$

By Lemma 36

$$m_n \doteq \inf_{s \in [0, T \wedge \tau_{n,L}]} f_n(s) \rightarrow +\infty \text{ as } n \rightarrow \infty. \quad (7.70)$$

By Ito's isometry, for  $i \leq k-1$

$$\begin{aligned} \mathbf{E} \sup_{s \in [0, T \wedge \tau_{n,L}]} \left| \int_0^s Z_{n,i}(s-) d(\sqrt{n}M_{n,i})(s) \right|^2 &\leq 4\mathbf{E} \int_0^{T \wedge \tau_{n,L}} Z_{n,i}^2(s-) d[\sqrt{n}M_{n,i}]_s \\ &\leq 4L^2 \mathbf{E}[\sqrt{n}M_{n,i}]_T = 4L^2 \mathbf{E} \langle \sqrt{n}M_{n,i} \rangle_T. \end{aligned}$$

where the second to last inequality is obtained by using  $\|\mathbf{Z}_n\|_{2,T\wedge\tau_{n,L}} \leq L$ . From the proof of Lemma 27 we see that for any  $i \leq k-1$ ,  $\mathbf{E} \langle \sqrt{n}M_{n,i} \rangle_T \rightarrow 0$  as  $n \rightarrow \infty$ . Hence from the definition of  $B_n$  and  $\vec{W}_n$

$$|B_n|_{*,T\wedge\tau_{n,L}} \xrightarrow{P} 0 \text{ as } n \rightarrow \infty. \quad (7.71)$$

Applying Lemma 33 to (7.69) with  $z(t) = \|X_n(t)\|^2$ ,  $f = 2f_n$ ,  $g = g_n$ ,  $M = B_n$ , and  $\tau = \tau_{n,L}$  shows for any  $t \in [0, T \wedge \tau_{n,L}]$

$$\|\vec{X}_n(t)\|^2 \leq \frac{|g_n|_{*,T\wedge\tau_{n,L}}}{2m_n} + 2|B_n|_{*,T\wedge\tau_{n,L}} + e^{-2m_n t} \left( \|\vec{X}_n(0)\|^2 + |B_n(0)| \right).$$

Taking  $t = \epsilon_n \doteq 1/\sqrt{m_n}$  and using (7.70), (7.71),  $|g_n|_{*,T\wedge\tau_{n,L}} \leq 2L^2$  and  $\vec{X}_n(0) \xrightarrow{P} (z_1, \dots, z_{k-1})^t$ , we see that

$$\sup_{t \in [\epsilon_n, T \wedge \tau_{n,L}]} \|\vec{X}_n(t)\| \xrightarrow{P} 0.$$

Since  $\epsilon_n \rightarrow 0$ , this shows (7.66) for any fixed  $\epsilon > 0$ . Finally, from (7.69), we see that

$$\sup_{t \in [0, \epsilon_n \wedge \tau_{n,L} \wedge T]} \|\vec{X}_n(t)\|^2 \leq \|\vec{X}_n(0)\|^2 + |g_n|_{*,T\wedge\tau_{n,L}} \epsilon_n + |B_n|_{*,T\wedge\tau_{n,L}}.$$

The convergence in (7.65) is now immediate on using that  $\epsilon_n \rightarrow 0$ ,  $|g_n|_{*,T\wedge\tau_{n,L}} \leq 2L^2$  and (7.71) holds.  $\square$

**Corollary 11.** *Under the assumptions of Lemma 37, for each  $i < k$ ,  $\int_0^{T \wedge \tau_{n,L}} |Z_{n,i}(s)| ds \xrightarrow{P} 0$ , as  $n \rightarrow \infty$ .*

*Proof.* For any  $\epsilon > 0$

$$\begin{aligned} \int_0^{T \wedge \tau_{n,L}} |Z_{n,i}(s)| ds &\leq \int_{[0, \epsilon \wedge \tau_{n,L}]} |Z_{n,i}(s)| ds + \int_{[\epsilon, T \wedge \tau_{n,L}]} |Z_{n,i}(s)| ds \\ &\leq L\epsilon + \sup_{s \in [\epsilon, T \wedge \tau_{n,L}]} |Z_{n,i}(s)| T. \end{aligned}$$

Now fix  $\delta > 0$  and let  $\epsilon = \frac{\delta}{2L}$ . Then for any  $i < k$

$$\mathbf{P} \left( \int_0^{T \wedge \tau_{n,L}} |Z_{n,i}(s)| ds > \delta \right) \leq \mathbf{P} \left( \sup_{s \in [\epsilon, T \wedge \tau_{n,L}]} |Z_{n,i}(s)| > \frac{\delta}{2T} \right), \quad (7.72)$$

which from (7.66) converges to 0 as  $n \rightarrow \infty$ . Since  $\delta > 0$  was arbitrary, this completes the proof.  $\square$

*Proof of Theorem 9.* Recall the conditions in the theorem. By Remark 7.2.3(ii) and the tightness of  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$ , the hypothesis of Lemma 27 holds. Hence by Skorokhod's embedding theorem, we can assume that  $\{(\mathbf{Z}_n(0), \mathbf{M}_n)\}_{n \in \mathbb{N}}$  and a standard Brownian motion  $B$  are defined on a common probability space such that for any  $T > 0$

$$\sup_{t \leq T} \left\| \sqrt{n} \mathbf{M}_n(t) - \sqrt{2} B(t) \mathbf{e}_k \right\|_2 \rightarrow 0 \quad (7.73)$$

and

$$\|\mathbf{Z}_n(0) - \mathbf{z}\|_2 \rightarrow 0 \quad (7.74)$$

almost surely, as  $n \rightarrow \infty$ . Let  $\mathbf{Y}$  and  $\mathbf{Y}_n$  be as in the statement of the theorem. Taking  $m \doteq r - k + 1$ ,  $\vec{Y}_n \doteq (\sum_{i=1}^k Z_{n,i}, Z_{n,k+1}, \dots, Z_{n,r})$  be the stochastic process with sample paths in  $\mathbb{D}([0, T] : \mathbb{R}^m)$  corresponding to the first  $m$  coordinates of  $\mathbf{Y}_n$ . Note  $\mathbf{Y}_{n,m+} = \mathbf{Z}_{n,r+}$ ,  $Z_{n,k} = Y_{n,1} - \sum_{i=1}^{k-1} Z_i$ , and for  $k = 1$ ,  $Y_{n,1} = Z_{n,1}$ . Hence by Corollary 10,  $\vec{Y}_n$  satisfy

$$\begin{aligned} Y_{n,1}(t) &= Y_{n,1}(0) - \int_0^t a_{n,k}(s) Y_{n,1}(s) ds - \mathbb{I}_{\{k=1\}} \int_0^t Y_{n,1}(s) ds + \int_0^t Y_{n,2}(s) ds + \sqrt{n} M_{n,k}(t) \\ &\quad + \sum_{i=1}^{k-1} \int_0^t a_{n,k}(s) Z_{n,i}(s) ds - \mathbb{I}_{\{k>1\}} \int_0^t Z_{n,1}(s) ds + \sum_{i=1}^{k-1} \sqrt{n} M_{n,i}(t), \end{aligned} \quad (7.75)$$

$$\begin{aligned} Y_{n,2}(t) &= Y_{n,2}(0) + \int_0^t a_{n,k}(s) Y_{n,1}(s) ds - \int_0^t Y_{n,2}(s) ds + \int_0^t Y_{n,3}(s) ds \\ &\quad - \sum_{i=1}^{k-1} \int_0^t a_{n,k}(s) Z_{n,i}(s) ds - \int_0^t \delta_{n,k+1}(s) Y_{n,2}(s) ds + \sqrt{n} M_{n,k+1}(t), \end{aligned} \quad (7.76)$$

and for  $i \in \{3, 4, \dots, m\}$

$$\begin{aligned} Y_{n,i}(t) &= Y_{n,i}(0) - \int_0^t Y_{n,i}(s) ds + \int_0^t Y_{n,i+1}(s) ds \\ &\quad + \int_0^t \delta_{n,k+i-2}(s) Y_{n,i-1}(s) ds - \int_0^t \delta_{n,k+i-1}(s) Y_{n,i}(s) ds + \sqrt{n} M_{n,k+i-1}(t). \end{aligned} \quad (7.77)$$

where  $a_{n,k}(s)$  is as in Lemma 36 and  $\delta_{n,i}(s) \doteq \Delta_{n,i}(Z_{n,i}(s))$  for  $i \in \mathbb{N}$ .

Since  $\|\mathbf{Z}_n\|_{2,T \wedge \tau_{n,L}} \leq L$ , we have by (7.60) that, for any  $i \in \mathbb{N}$ ,

$$|\delta_{n,i}|_{*,T \wedge \tau_{n,L}} \leq \gamma_{n,L} \rightarrow 0 \text{ as } n \rightarrow \infty. \quad (7.78)$$

Moreover since  $\beta'_n(\mu_{n,k}) \rightarrow \alpha \in [0, \infty)$ , this also shows that the term

$$\sup_{s \in [0, T \wedge \tau_{n,L}]} |a_{n,k}(s) - \alpha| \rightarrow 0 \text{ as } n \rightarrow \infty. \quad (7.79)$$

We now show that

$$\|\mathbf{Y}_n - \mathbf{Y}\|_{2, T \wedge \tau_{n,L}} \xrightarrow{P} 0 \text{ as } n \rightarrow \infty. \quad (7.80)$$

To see this, note that, by Remark 7.2.3(ii), the hypothesis of Lemma 29 is satisfied, and hence

$\|\mathbf{Z}_{n,r+}\|_{2, T \wedge \tau_{n,L}} \xrightarrow{P} 0$ . Since  $\mathbf{Y}_{n,m+} = \mathbf{Z}_{n,r+}$  and  $\mathbf{Y}_{m+} = 0$ , this shows that

$$\|\mathbf{Y}_{n,m+} - \mathbf{Y}_{m+}\|_{2, T \wedge \tau_{n,L}} \xrightarrow{P} 0. \quad (7.81)$$

Thus in order to prove (7.80) it suffices to show that  $\sum_{i=1}^m |Y_{n,i} - Y_i|_{*, T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . To show this we consider  $U_{n,i} \doteq Y_{n,i} - Y_i$ . Subtracting (7.10) from (7.75), (7.76) and (7.77), we see

$$\begin{aligned} U_{n,1}(t) &= U_{n,1}(0) - (\alpha + \mathbb{I}_{\{k=1\}}) \int_0^t U_{n,1}(s) ds + \int_0^t U_{n,2}(s) ds + \sqrt{n} M_{n,k}(t) - B(t) + W_{n,1}(t) \\ U_{n,2}(t) &= U_{n,2}(0) + \alpha \int_0^t U_{n,1}(s) ds - \int_0^t U_{n,2}(s) ds + \int_0^t U_{n,3}(s) ds + W_{n,2}(t) \\ U_{n,i}(t) &= U_{n,i}(0) - \int_0^t U_{n,i}(s) ds + \int_0^t U_{n,i+1}(s) ds + W_{n,i}(t) \quad \text{for } i \in \{3, 4, \dots, m\} \end{aligned} \quad (7.82)$$

where

$$\begin{aligned} W_{n,1}(t) &\doteq \int_0^t (\alpha - a_{n,k}(s)) Y_{n,1}(s) ds + \sum_{i=1}^{k-1} \int_0^t a_{n,k}(s) Z_{n,i}(s) ds - \mathbb{I}_{\{k>1\}} \int_0^t Z_{n,1}(s) ds + \sum_{i=1}^{k-1} \sqrt{n} M_{n,i}(t) \\ W_{n,2}(t) &\doteq \int_0^t (a_{n,k}(s) - \alpha) Y_{n,1}(s) ds - \sum_{i=1}^{k-1} \int_0^t a_{n,k}(s) Z_{n,i}(s) ds - \int_0^t \delta_{n,k+1}(s) Y_{n,2}(s) ds + \sqrt{n} M_{n,k+1}(t), \\ W_{n,i}(t) &\doteq \int_0^t \delta_{n,k+i-2}(s) Y_{n,i-1}(s) ds - \int_0^t \delta_{n,k+i-1}(s) Y_{n,i}(s) ds + \sqrt{n} M_{n,k+i-1}(t) \quad \text{for } i \in \{3, \dots, m\}. \end{aligned}$$

Note that, for each  $n$ ,  $\|\mathbf{Y}_n\|_{2, T \wedge \tau_{n,L}} \leq k \|\mathbf{Z}_n\|_{2, T \wedge \tau_{n,L}}$ , which by (7.51) is bounded above by  $kL$ .

Hence by (7.79), (7.78), (7.73) and Corollary 11,

$$|W_{n,i}|_{*, T \wedge \tau_{n,L}} \xrightarrow{P} 0 \text{ as } n \rightarrow \infty \quad (7.83)$$



for each  $i \in [m]$ . Let  $\|U_n\|_{1,t} \doteq \sup_{s \in [0,t]} \sum_{i=1}^m |U_{n,i}(s)|$ . Then, from (7.82), for any  $t \in [0, T \wedge \tau_{n,L}]$

$$\|U_n\|_{1,t} \leq \sum_{i=1}^m \left( |U_{n,i}(0)| + |W_{n,i}|_{*,T \wedge \tau_{n,L}} \right) + |\sqrt{n}M_{n,k} - B|_{*,T} + R \int_0^t \|U_n\|_{1,s} ds$$

with  $R \doteq \max(2\alpha + \mathbb{I}_{\{k=1\}}, 2)$ . Hence by Gronwall's inequality

$$\|U_n\|_{1,T \wedge \tau_{n,L}} \leq \left( |\sqrt{n}M_{n,k} - B|_{*,T} + \sum_{i=1}^m \left( |U_{n,i}(0)| + |W_{n,i}|_{*,T \wedge \tau_{n,L}} \right) \right) e^{RT}.$$

By our hypothesis, as  $n \rightarrow \infty$ ,  $|U_{n,i}(0)| = |Z_{n,k+i-1}(0) - z_{n,k+i-1}| \xrightarrow{P} 0$  for each  $i \in [m]$ . Hence by (7.83) and (7.73),  $\|U_n\|_{1,T \wedge \tau_{n,L}} = \sum_{i=1}^m |Y_{n,i} - Y_i|_{*,T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . Combined with (7.81), this completes the proof of (7.80).

Next we prove (7.9). Fix  $\delta > 0$ . Since  $\mathbf{Y}$  has sample paths in  $C([0, T] : l_2)$ , we can find  $L_1 \in (0, \infty)$  so that

$$\mathbf{P}\left(\|\mathbf{Y}\|_{2,T} > L_1\right) \leq \frac{\delta}{2}. \quad (7.84)$$

Also, since  $\mathbf{Z}_n(0) \xrightarrow{P} \mathbf{z}$ , we can find a  $L_2 \in (0, \infty)$  so that

$$\sup_n \mathbf{P}(\|\mathbf{Z}_n(0)\|_2 > L_2) \leq \frac{\delta}{2}. \quad (7.85)$$

Let  $L \doteq (L_1 + 1) + k(L_2 + 1) + 1$ . Also, let  $\vec{X}_n$  be as in Lemma 37 when  $k > 1$ . For  $k = 1$ , we set  $\vec{X}_n \doteq 0$ . Then,

$$\begin{aligned} \|\mathbf{Z}_n\|_{2,T \wedge \tau_{n,L}} &\leq \left\| \vec{X}_n \right\|_{2,T \wedge \tau_{n,L}} + \left\| \mathbf{Y}_n - \mathbf{e}_1 \sum_{i=1}^{k-1} Z_{n,i} \right\|_{2,T \wedge \tau_{n,L}} \\ &\leq \|\mathbf{Y}_n\|_{2,T \wedge \tau_{n,L}} + k \mathbb{I}_{\{k>1\}} \left\| \vec{X}_n \right\|_{2,T \wedge \tau_{n,L}}. \end{aligned}$$

Hence for each  $n \in \mathbb{N}$

$$\begin{aligned} \mathbf{P}(\tau_{n,L} \leq T) &\leq \mathbf{P}\left(\|\mathbf{Z}_n\|_{2,T \wedge \tau_{n,L}} > L - 1\right) \\ &\leq \mathbf{P}\left(\|\mathbf{Y}_n\|_{2,T \wedge \tau_{n,L}} > L_1 + 1\right) + \mathbf{P}\left(\left\| \vec{X}_n \right\|_{2,T \wedge \tau_{n,L}} > L_2 + 1\right), \\ &\leq \delta + \mathbf{P}\left(\|\mathbf{Y}_n - \mathbf{Y}\|_{2,T \wedge \tau_{n,L}} > 1\right) + \mathbf{P}\left(\left\| \vec{X}_n \right\|_{2,T \wedge \tau_{n,L}} > \left\| \vec{X}_n(0) \right\| + 1\right), \end{aligned}$$

where the last inequality uses (7.84) and (7.85). From Lemma 37 and (7.80) we see

$$\limsup_{n \rightarrow \infty} \mathbf{P}\left(\|\mathbf{Z}_n\|_{2,T} \geq L\right) \leq \limsup_{n \rightarrow \infty} \mathbf{P}(\tau_{n,L} \leq T) \leq \delta.$$

Since  $\delta > 0$  is arbitrary, the convergence in (7.9) is now immediate.

This convergence in particular says that  $\lim_{L \rightarrow \infty} \sup_n \mathbf{P}(\tau_{n,L} \leq T) = 0$ . Using Lemma 30 with  $F_n(t) = \|\mathbf{Y}_n - \mathbf{Y}\|_{2,t}$  we now see from (7.80) that  $\|\mathbf{Y}_n - \mathbf{Y}\|_{2,T} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . Similarly, if  $k > 1$ , then taking  $F_n(t) = \sup_{s \in [\epsilon, t]} |Z_{n,i}(s)|$  in Lemma 30 we conclude from Lemma 37 that for each  $i \in [k-1]$  and  $\epsilon > 0$   $\sup_{s \in [\epsilon, T]} |Z_{n,i}(s)| \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . This completes the proof of Theorem 9.  $\square$

## 7.8 Proof of Theorem 10

In this section we give the proof of Theorem 10. We begin by giving a convenient representation for  $\mathbf{Z}_n$  under the assumptions of Theorem 10 and establishing some apriori convergence properties.

**Lemma 38.** *Suppose  $c_n = \frac{d_n}{\sqrt{n}} \rightarrow c \in (0, \infty)$  and  $\lambda_n = 1 - \left(\frac{\log d_n}{d_n} + \frac{\alpha_n}{\sqrt{n}}\right)$  where  $\alpha_n \in \mathbb{R}$ ,  $\liminf_{n \rightarrow \infty} \alpha_n > -\infty$  and  $\frac{\alpha_n}{n^{1/4}} \rightarrow 0$ . Suppose also that  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$  is a tight sequence of random variables and  $\mathbf{Z}_{n,r+}(0) \xrightarrow{P} \mathbf{0}$  in  $l_2$  for some  $r \geq 2$ . Then there are real stochastic processes  $\delta_n, W_{n,i}$  with sample paths in  $\mathbb{D}([0, \infty) : \mathbb{R})$  such that for any  $t \geq 0$*

$$\begin{aligned} Z_{n,1}(t) &= Z_{n,1}(0) - \int_0^t Z_{n,1}(s) ds + \int_0^t Z_{n,2}(s) ds + \sqrt{n} M_{n,1}(t) \\ &\quad - (c_n e^{c_n \alpha_n})^{-1} \int_0^t (1 + \delta_n(s)) (e^{c_n Z_{n,1}(s)} - 1) ds \\ Z_{n,2}(t) &= Z_{n,2}(0) - \int_0^t Z_{n,2}(s) ds + \int_0^t Z_{n,3}(s) ds + W_{n,2}(t) \\ &\quad + (c_n e^{c_n \alpha_n})^{-1} \int_0^t (1 + \delta_n(s)) (e^{c_n Z_{n,1}(s)} - 1) ds \\ Z_{n,i}(t) &= Z_{n,i}(0) - \int_0^t Z_{n,i}(s) ds + \int_0^t Z_{n,i+1}(s) ds + W_{n,i}(t) \quad \text{for } i \in \{3, \dots, r\} \end{aligned} \tag{7.86}$$

and for any fixed  $L, T \in (0, \infty)$ ,

1.  $\sqrt{n} M_{n,1} \Rightarrow \sqrt{2} B$  in  $\mathbb{D}([0, \infty) : \mathbb{R})$  where  $B$  is a standard Brownian motion,
2.  $|\delta_n|_{*, T_n} \rightarrow 0$  a.s.

$$3. |W_{n,i}|_{*,T_n} \xrightarrow{P} 0 \text{ for } i \in \{2, \dots, r\},$$

$$4. \|\mathbf{Z}_{n,r+}\|_{2,T_n} \xrightarrow{P} 0,$$

where  $T_n \doteq T \wedge \tau_{n,L}$  and  $\tau_{n,L}$  is defined as in (7.50).

*Proof.* Recall the definition of  $t_{n,i}$  from Lemma 26. Define

$$\delta_n(s) \doteq t_{n,1}(Z_{n,1}(s))c_n \left( e^{c_n[Z_{n,1}(s)-\alpha_n]} - e^{-c_n\alpha_n} \right)^{-1} - 1$$

so that

$$t_{n,1}(Z_{n,1}(s)) = (1 + \delta_n(s))c_n^{-1} \left( e^{c_n[Z_{n,1}(s)-\alpha_n]} - e^{-c_n\alpha_n} \right).$$

Since  $\sup_{s \leq T \wedge \tau_{n,L}} |Z_{n,1}(s)| \leq L$ , Lemma 32 shows that  $|\delta_n|_{*,T_n} \rightarrow 0$  a.s. Define

$$\begin{aligned} W_{n,2}(t) &\doteq - \int_0^t t_{n,2}(Z_{n,2}(s))ds + \sqrt{n}M_{n,2}(t) \\ W_{n,i}(t) &\doteq \int_0^t t_{n,i-1}(Z_{n,i-1}(s))ds - \int_0^t t_{n,i}(Z_{n,i}(s))ds + \sqrt{n}M_{n,i}(t) \quad \text{for } i \in \{3, \dots, r\}. \end{aligned}$$

From Lemma 26 it follows that (7.86) is satisfied. Lemma 31 shows that  $\boldsymbol{\mu}_n \rightarrow \mathbf{f}_1 \in l_1^\downarrow$ . Along with the assumed tightness of  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$ , this shows  $\mathbf{G}_n(0) = \boldsymbol{\mu}_n + \frac{\mathbf{Z}_n(0)}{\sqrt{n}} \xrightarrow{P} \mathbf{f}_1$  in  $l_1^\downarrow$ . Hence by Lemma 27 and Lemma 29,

$$\sqrt{n}\mathbf{M}_n \Rightarrow \sqrt{2}B\mathbf{e}_1 \text{ in } \mathbb{D}([0, \infty) : l_2) \quad (7.87)$$

and  $\|\mathbf{Z}_{n,r+}\|_{2,T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . Since  $|Z_{n,i}|_{*,T \wedge \tau_{n,L}} \leq L$  and  $\mu_{n,2} \rightarrow 0$ , Lemma 28, together with (7.87), shows that  $|W_{n,i}|_{*,T_n} \xrightarrow{P} 0$  for each  $i \in \{2, \dots, r\}$ , as  $n \rightarrow \infty$ .  $\square$

The next lemma gives pathwise existence and uniqueness of solutions to a system of stochastic differential equations in which the drift fails to satisfy a linear growth condition.

**Lemma 39.** Suppose  $c \in (0, \infty)$ ,  $\alpha \in (0, \infty]$  and  $B$  is a standard Brownian motion. Then for any  $r \geq 2$  the system of equations

$$\begin{aligned} Z_1(t) &= z_1 - \int_0^t Z_1(s)ds + \int_0^t Z_2(s)ds + \sqrt{2}B(t) - (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds \\ Z_2(t) &= z_2 - \int_0^t Z_2(s)ds + \int_0^t Z_3(s)ds + (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds \\ Z_i(t) &= z_i - \int_0^t Z_i(s)ds + \int_0^t Z_{i+1}(s)ds \quad \text{for } i \in \{3, \dots, r\} \\ Z_i(t) &= 0 \quad \text{for } i > r \end{aligned} \tag{7.88}$$

has a unique pathwise solution  $\mathbf{Z}$  with sample paths in  $C([0, \infty) : l_2)$  for any  $(z_1, \dots, z_r) \in \mathbb{R}^r$ .

*Proof.* The case when  $\alpha = \infty$  is standard and is thus omitted. Consider now the case  $\alpha < \infty$ . It is straightforward to see that there is a unique  $\mathbf{Z}_{2+} \doteq (Z_3, Z_4, \dots)$  in  $C([0, \infty) : l_2)$  that solves the last two equations in (7.88). Hence it suffices to show that, the system of equations

$$\begin{aligned} Z_1(t) &= z_1 - (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds + \int_0^t (Z_2(s) - Z_1(s))ds + \sqrt{2}B(t) \\ Z_2(t) &= z_2 + (ce^{c\alpha})^{-1} \int_0^t (e^{cZ_1(s)} - 1)ds - \int_0^t Z_2(s)ds + \int_0^t f(s)ds \end{aligned} \tag{7.89}$$

has a unique pathwise solution  $(Z_1, Z_2)$  with sample paths in  $C([0, \infty) : \mathbb{R}^2)$  where  $f \doteq Z_3 \in C([0, \infty) : \mathbb{R})$  is a given (non-random) continuous trajectory and  $(z_1, z_2) \in \mathbb{R}^2$ .

Define  $y_1 = z_1 + z_2$ ,  $y_2 = z_2$  and consider the equation:

$$\begin{aligned} Y_1(t) &= y_1 - (ce^{c\alpha})^{-1} \int_0^t (e^{cY_1(s)} - 1)ds + \int_0^t (Y_2(s) - 2Y_1(s))ds + \sqrt{2}B(t) \\ Y_2(t) &= y_2 - \int_0^t Y_1(s)ds + \int_0^t f(s)ds + \sqrt{2}B(t). \end{aligned} \tag{7.90}$$

Note that  $(Z_1, Z_2)$  solve (7.89) if and only if  $(Y_1, Y_2)$ , with  $Y_1 = Z_1$  and  $Y_2 = Z_1 + Z_2$  solve (7.90). Thus it suffices to prove existence and uniqueness of solutions for (7.90).

For  $L \in (0, \infty)$ , let  $\eta_L : \mathbb{R} \rightarrow [0, 1]$  be such that  $\eta_L$  is smooth,  $\eta_L(x) = 1$  for  $|x| \leq L$  and  $\eta_L(x) = 0$  for  $|x| \geq L + 1$ . Consider the equation

$$\begin{aligned} Y_1^L(t) &= y_1 - (ce^{c\alpha})^{-1} \int_0^t e^{cY_1^L(s)} \eta_L(Y_1^L(s)) ds + (ce^{c\alpha})^{-1} t \\ &\quad + \int_0^t (Y_2^L(s) - 2Y_1^L(s)) ds + \sqrt{2}B(t) \\ Y_2^L(t) &= y_2 - \int_0^t Y_1^L(s) ds + \int_0^t f(s) ds + \sqrt{2}B(t). \end{aligned} \tag{7.91}$$

Since for each  $L$  (7.91) is an equation with (globally) Lipschitz coefficients, by standard results, it has a unique pathwise continuous solution.

Fix  $T \in (0, \infty)$  and let  $\tau_L = \inf\{t \geq 0 : |Y_1^L(t)| \geq L\} \wedge T$  for any  $L > 0$ . Then by pathwise uniqueness of (7.91), for  $0 \leq t \leq \tau_L \wedge \tau_{L+1}$ ,

$$Y^L(t) = Y^{L+1}(t).$$

This in particular shows that,  $\tau_L \leq \tau_{L+1}$  a.s.

We now estimate the second moment of  $|Y_1^L(t)|$ . By Itô's formula

$$\begin{aligned} (Y_1^L(t))^2 &= (y_1)^2 - 2(ce^{c\alpha})^{-1} \int_0^t Y_1^L(s) e^{cY_1^L(s)} \eta_L(Y_1^L(s)) ds + 2(ce^{c\alpha})^{-1} \int_0^t Y_1^L(s) ds \\ &\quad + 2 \int_0^t Y_1^L(s) (Y_2^L(s) - 2Y_1^L(s)) ds + 2\sqrt{2} \int_0^t Y_1^L(s) dB(s) + 2t \\ (Y_2^L(t))^2 &= (y_2)^2 - 2 \int_0^t Y_1^L(s) Y_2^L(s) ds + 2 \int_0^t Y_2^L(s) f(s) ds + 2\sqrt{2} \int_0^t Y_2^L(s) dB(s) + 2t. \end{aligned}$$

Thus

$$\begin{aligned} (Y_1^L(t))^2 + (Y_2^L(t))^2 &= (y_1)^2 + (y_2)^2 - 2(ce^{c\alpha})^{-1} \int_0^t Y_1^L(s) e^{cY_1^L(s)} \eta_L(Y_1^L(s)) ds \\ &\quad + 2(ce^{c\alpha})^{-1} \int_0^t Y_1^L(s) ds + 2 \int_0^t Y_2^L(s) f(s) ds \\ &\quad - 4 \int_0^t (Y_1^L(s))^2 ds + 2\sqrt{2} \int_0^t (Y_1^L(s) + Y_2^L(s)) dB(s) + 4t. \end{aligned}$$

Since  $c > 0$ , we have on using the inequality  $|x| \leq 1 + |x|^2$  that  $-xe^{cx}\eta_L(x) \leq (1 + |x|^2)$  for all  $x \in \mathbb{R}$ . Thus with  $\|Y^L\|_{*,t} \doteq \sup_{s \in [0,t]} \|Y^L(s)\|$ :

$$\begin{aligned} \|Y^L\|_{*,t}^2 &\leq \|y\|^2 + 4(ce^{c\alpha})^{-1} \int_0^t (1 + \|Y^L\|_{*,s}^2) ds \\ &\quad + 2 \int_0^t (1 + \|Y^L\|_{*,s}^2) |f(s)| ds \\ &\quad + 2\sqrt{2} \left( 1 + \sup_{0 \leq s \leq t} \left| \int_0^s (Y_1^L(u) + Y_2^L(u)) dB(u) \right|^2 \right) + 4t. \end{aligned}$$

Taking expectations, for any  $t \in [0, T]$

$$\begin{aligned} \mathbf{E}\|Y^L\|_{*,t}^2 &\leq \|y\|^2 + (4(ce^{c\alpha})^{-1} + 2|f|_{*,T}) \int_0^t (1 + \mathbf{E}\|Y^L\|_{*,s}^2) ds \\ &\quad + 2\sqrt{2} \left( 1 + 4\mathbf{E} \int_0^t |Y_1^L(u) + Y_2^L(u)|^2 du \right) + 4t \\ &\leq (\|y\|^2 + K(T+1)) + K \int_0^t \mathbf{E}\|Y^L\|_{*,s}^2 ds. \end{aligned}$$

with  $K \doteq 4(ce^{c\alpha})^{-1} + 2|f|_{*,T} + 16\sqrt{2}$ . By Gronwall lemma, for every  $L \in \mathbb{N}$

$$\mathbf{E}\|Y^L\|_{*,T}^2 \leq (\|y\|^2 + K(T+1))e^{KT} \doteq c_1.$$

Thus, as  $L \rightarrow \infty$

$$P(\tau_L < T) \leq P(\|Y^L\|_{*,T} \geq L) \leq c_1/L^2 \rightarrow 0$$

and consequently  $\tau_L \uparrow T$  as  $L \rightarrow \infty$ . Now define  $Y(t) \doteq Y^L(t)$  for  $0 \leq t \leq \tau_L$ . Then  $Y$  is a solution of (7.90) on  $[0, T]$ . The same argument as before shows that this is the unique pathwise solution on  $[0, T]$ . Since  $T$  is arbitrary we get a unique pathwise solution of (7.90) on  $[0, \infty)$ . This completes the proof of the lemma.  $\square$

**Lemma 40.** *Suppose the assumptions of Theorem 10 hold. Suppose further that  $\mathbf{Z}_n(0), \mathbf{M}_n$  and a standard Brownian motion  $B$  are given on a common probability space such that  $\mathbf{Z}_n(0) \rightarrow \mathbf{z}$  in  $l_1^1$  and  $\mathbf{M}_n \rightarrow \sqrt{2}B\mathbf{e}_1$  in  $\mathbb{D}([0, \infty) : l_2)$  almost surely. Let  $\mathbf{Z}$  be as defined in Lemma 39. Then for any  $T, L \in (0, \infty)$*

$$\|\mathbf{Z}_n - \mathbf{Z}\|_{2, T \wedge \tau_{n,L} \wedge \tau_L} \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty. \quad (7.92)$$

where  $\tau_L \doteq \inf \left\{ t \mid \|\mathbf{Z}_n(t)\|_{2,t} > L \right\}$ .

*Proof.* Fix  $L, T \in (0, \infty)$  and let  $T_n \doteq T \wedge \tau_{n,L} \wedge \tau_L$  and  $U_{n,i} \doteq Z_{n,i} - Z_i$  for  $i \in \mathbb{N}$ . Using the estimate  $|e^{ax} - e^{ay}| \leq ae^{a(x \vee y)} |x - y|$  for  $x, y \in \mathbb{R}$ ,  $a \geq 0$  and since  $|Z_{n,1}(s)|, |Z_1(s)| \leq L$  for any  $s \in [0, T_n]$ ,

$$\begin{aligned} & \left| a_n(s) e^{c_n Z_{n,1}(s)} - a e^{c Z_1(s)} \right| \\ & \leq \left| a_n(s) e^{c_n Z_{n,1}(s)} - a_n(s) e^{c_n Z_1(s)} \right| + \left| a_n(s) e^{c_n Z_1(s)} - a_n(s) e^{c Z_1(s)} \right| + \left| e^{c Z_1(s)} \right| |a_n(s) - a| \\ & \leq |a_n(s)| c_n e^{c_n L} |U_{n,1}(s)| + |a_n(s)| L e^{L(c_n \vee c)} |c_n - c| + e^{cL} |a_n(s) - a| \end{aligned}$$

where  $a_n(s) \doteq (c_n e^{c_n \alpha_n})^{-1} (1 + \delta_n(s))$ ,  $c_n = d_n / \sqrt{n}$ ,  $\delta_n$  is as in Lemma 38, and  $a \doteq (c e^{c\alpha})^{-1}$ . Since  $c_n \rightarrow c$  and  $|\delta_n|_{*,T_n} \rightarrow 0$  by Lemma 38,  $|a_n - a|_{*,T_n} \rightarrow 0$ . Hence for any  $s \in [0, T_n]$

$$|a_n e^{c_n Z_{n,1}} - a e^{c Z_1}|_{*,s} \leq K |U_{n,1}|_{*,s} + r_n \quad (7.93)$$

where  $K \doteq \sup_n (c_n e^{c_n L} |a_n|_{*,T_n}) < \infty$  and

$$r_n \doteq |a_n|_{*,T_n} L e^{L(c_n \vee c)} |c_n - c| + e^{cL} |a_n - a|_{*,T_n} \rightarrow 0$$

almost surely. Subtracting (7.88) from (7.86), for any  $t > 0$ ,

$$\begin{aligned} U_{n,1}(t) &= U_{n,1}(0) - \int_0^t (U_{n,1}(s) - U_{n,2}(s)) ds + \sqrt{n} M_{n,1}(t) - \sqrt{2} B(t) \\ &\quad - \int_0^t \left( a_{n,1}(s) e^{c_n Z_{n,1}(s)} - a e^{c Z_1(s)} \right) ds + \int_0^t (a_n(s) - a) ds \\ U_{n,2}(t) &= U_{n,2}(0) - \int_0^t (U_{n,2}(s) - U_{n,3}(s)) ds + W_{n,2}(t) \\ &\quad + \int_0^t \left( a_{n,1}(s) e^{c_n Z_{n,1}(s)} - a e^{c Z_1(s)} \right) ds - \int_0^t (a_n(s) - a) ds \\ U_{n,i}(t) &= U_{n,i}(0) - \int_0^t (U_{n,i}(s) - U_{n,i+1}(s)) ds + W_{n,i}(t) \quad \text{for } i \in \{3, \dots, r\}. \end{aligned} \quad (7.94)$$

Let  $H_t \doteq \sup_{s \in [0,t]} \sum_{i=1}^r |U_{n,i}(s)|$ . Then from (7.93) and (7.94), for any  $t \in [0, T_n]$ ,

$$H_t \leq H_0 + \left| \sqrt{n} M_{n,1} - \sqrt{2} B \right|_{*,T} + 2T(|a_n - a|_{*,T_n} + r_n) + \sum_{i=2}^r |W_{n,i}|_{*,T_n} + |U_{n,r+1}|_{*,T_n} + 2(1+K) \int_0^t H_s ds.$$

Hence by Gronwall's lemma

$$H_{T_n} \leq \left( H_0 + \left| \sqrt{n}M_{n,1} - \sqrt{2}B \right|_{*,T} + 2T(|a_n - a|_{*,T_n} + r_n) + \sum_{i=2}^r |W_{n,i}|_{*,T_n} + |U_{n,r+1}|_{*,T_n} \right) e^{2(1+K)T}.$$

Note  $\mathbf{U}_{n,r+} = \mathbf{Z}_{n,r+}$  and  $U_{n,i}(0) = Z_{n,i}(0) - z_i$ . Then by Lemma 38,  $H_{T_n} \xrightarrow{P} 0$  and  $\|\mathbf{U}_{n,r+}\|_{2,T_n} \xrightarrow{P} 0$ . Together these show  $\|\mathbf{U}\|_{2,T_n} = \|\mathbf{Z}_n - \mathbf{Z}\|_{2,T_n} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ .  $\square$

**Corollary 12.** *Under assumptions of Theorem 10,  $\left\{ \|\mathbf{Z}_n\|_{2,T} \right\}_{n \in \mathbb{N}}$  is a tight sequence of random variables and*

$$\lim_{L \rightarrow \infty} \sup_n \mathbf{P}(\tau_{n,L} \leq T) = 0. \quad (7.95)$$

*Proof.* Fix  $\delta > 0$ . Since  $\mathbf{Z}$  has sample paths in  $C([0, T] : l_2)$ , we can find  $L \in (0, \infty)$  so that  $\mathbf{P}(\|\mathbf{Z}\|_{2,T} > L) \leq \delta$ . Note that the events  $\left\{ \|\mathbf{Z}\|_{2,T} \leq L \right\} \subseteq \{\tau_{L+2} > T\}$ , where  $\tau_{L+2} \doteq \inf\{t \mid \|\mathbf{Z}(t)\|_2 > L + 2\}$ . Hence for each  $n \in \mathbb{N}$ , by right continuity of  $\mathbf{Z}_n$

$$\begin{aligned} \mathbf{P}(\tau_{n,L+2} \leq T) &\leq \mathbf{P}\left(\|\mathbf{Z}_n\|_{2,T \wedge \tau_{n,L+2}} > L + 1\right) \leq \mathbf{P}\left(\|\mathbf{Z}_n - \mathbf{Z}\|_{2,T \wedge \tau_{n,L+2}} > 1 \text{ or } \|\mathbf{Z}\|_{2,T} > L\right) \\ &\leq \mathbf{P}\left(\|\mathbf{Z}_n - \mathbf{Z}\|_{2,T \wedge \tau_{n,L+2} \wedge \tau_{L+2}} > 1\right) + \mathbf{P}\left(\|\mathbf{Z}\|_{2,T} > L\right) \\ &\leq \delta + \mathbf{P}\left(\|\mathbf{Z}_n - \mathbf{Z}\|_{2,T \wedge \tau_{n,L+2} \wedge \tau_{L+2}} > 1\right). \end{aligned}$$

Sending  $n \rightarrow \infty$  and using Lemma 40 we see that  $\limsup_n \mathbf{P}(\tau_{n,L+2} \leq T) \leq \delta$ . Finally,

$$\limsup_n \mathbf{P}\left(\|\mathbf{Z}_n\|_{2,T} > L + 2\right) \leq \limsup_n \mathbf{P}(\tau_{n,L+2} \leq T) \leq \delta.$$

Since  $\delta > 0$  is arbitrary, this shows that  $\left\{ \|\mathbf{Z}_n\|_{2,T} \right\}_{n \in \mathbb{N}}$  is tight. The convergence in (7.95) is now immediate.  $\square$

*Proof of Theorem 10.* Using Lemma 38 and Skorohod embedding theorem we can assume without loss of generality that  $\mathbf{Z}_n(0), \mathbf{M}_n$  and  $B$  are given on a common probability space,  $\mathbf{Z}_n(0) \rightarrow \mathbf{z}$  in  $l_1^\downarrow$ , and  $\mathbf{M}_n \rightarrow \sqrt{2}B\mathbf{e}_1$  in  $\mathbb{D}([0, \infty) : l_2)$  almost surely. From Lemma 40 we now have that for every  $T, L \in (0, \infty)$  (7.92) holds. Finally using from Corollary 12 the fact that  $\sup_n \mathbf{P}(\tau_{n,L} \leq T) + \mathbf{P}(\tau_L \leq T) \rightarrow 0$  as  $L \rightarrow \infty$ , we have that  $\|\mathbf{Z}_n - \mathbf{Z}\|_{2,T} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . The result follows.  $\square$



## 7.9 Proof of Theorem 11

In this section we give the proof of Theorem 11. As for the proof of Theorem 10 we begin with a convenient representation for  $\mathbf{Z}_n$  and by establishing some useful convergence properties.

**Lemma 41.** *Let  $\lambda_n, \alpha_n, d_n$  be as in the statement of Theorem 11. Suppose that  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$  is a tight sequence of random variables and  $\mathbf{Z}_{n,r+}(0) \xrightarrow{P} \mathbf{0}$  in  $l_2$  for some  $r \geq 2$ . Then there are real stochastic processes  $U_n, V_n, W_{n,i}, \eta_n$  with sample paths in  $\mathbb{D}([0, \infty) : \mathbb{R})$  so that,  $U_n, \eta_n$  have absolutely continuous paths a.s.,  $U_n(0) = \eta_n(0) = 0$ , and for any  $t \geq 0$*

$$Z_{n,1}(t) = Z_{n,1}(0) - \int_0^t Z_{n,1}(s)ds + \int_0^t Z_{n,2}(s)ds + \sqrt{n}M_{n,1}(t) + U_n(t) - \eta_n(t) \quad (7.96)$$

$$Z_{n,2}(t) = Z_{n,2}(0) - \int_0^t Z_{n,2}(s)ds + \int_0^t Z_{n,3}(s)ds + V_n(t) + \eta_n(t) \quad (7.97)$$

$$Z_{n,i}(t) = Z_{n,i}(0) - \int_0^t Z_{n,i}(s)ds + \int_0^t Z_{n,i+1}(s)ds + W_{n,i}(t) \quad \text{for } i \in \{3, \dots, r\}. \quad (7.98)$$

Furthermore,  $\eta_n$  is non-decreasing process with  $\eta_n(0) = 0$  that satisfies

$$\eta_n(t) = \int_0^t \mathbb{I}_{\{Z_{n,1}(s) \geq \theta_n\}} d\eta_n(s) \quad \text{a.s.} \quad (7.99)$$

for some constants  $\theta_n = \alpha_n + O(\sqrt{n}/d_n) \geq 0$  as  $n \rightarrow \infty$ . Also for any  $L, T \in (0, \infty)$ , as  $n \rightarrow \infty$

1.  $\sqrt{n}M_{n,1} \Rightarrow \sqrt{2}B$  in  $\mathbb{D}([0, T] : \mathbb{R})$
2.  $\text{TV}(U_n; [0, T_n]) \doteq \int_0^{T_n} |\dot{U}_n(s)|ds \xrightarrow{P} 0$
3.  $|V_n|_{*, T_n} \xrightarrow{P} 0$  and  $|W_{n,i}|_{*, T_n} \xrightarrow{P} 0$  for  $i \in \{3, \dots, r\}$
4.  $\|\mathbf{Z}_{n,r+}\|_{2, T_n} \xrightarrow{P} 0$ .

Here  $B$  is a standard Brownian motion and  $T_n \doteq T \wedge \tau_{n,L}$ .

*Proof.* By our assumptions on  $\alpha_n$ , we can find  $\kappa \in (0, \infty)$  be such that  $\theta_n \doteq \alpha_n + \frac{\kappa\sqrt{n}}{d_n} \geq 0$  for every  $n$ . Note  $\theta_n \rightarrow \alpha$  as  $n \rightarrow \infty$ . Define

$$U_n(t) \doteq - \int_0^t t_{n,1}(Z_{n,1}(s)) \mathbb{I}_{\{Z_{n,1}(s) < \theta_n\}} ds$$

$$\eta_n(t) \doteq \int_0^t t_{n,1}(Z_{n,1}(s)) \mathbb{I}_{\{Z_{n,1}(s) \geq \theta_n\}} ds$$

so that  $\eta_n(t) = \int_0^t \mathbb{I}_{\{Z_{n,1}(s) \geq \theta_n\}} d\eta_n(s)$ , and

$$\int_0^t t_{n,1}(Z_{n,1}(s)) ds = \eta_n(t) - U_n(t). \quad (7.100)$$

From Lemma 26 it then follows that (7.96) is satisfied. Recall the expression for  $t_{n,1}(z)$  from (7.46). Then, by monotonicity of  $\beta_n$ ,  $t_{n,1}(z) \geq 0$  whenever  $z \geq 0$ . Since  $\theta_n \geq 0$ ,  $\eta_n$  is non-decreasing and

$$\begin{aligned} \sup_{z \leq \theta_n} |t_{n,1}(z)| &\leq 2\sqrt{n}\beta_n(\lambda_n + \theta_n/\sqrt{n}) \leq 2\sqrt{n}(\lambda_n + \theta_n/\sqrt{n})^{d_n} \\ &= 2\sqrt{n}(1 - ((\log d_n)/d_n + (\alpha_n - \theta_n)/\sqrt{n}))^{d_n} = 2\sqrt{n}(1 - (\log d_n - \kappa)/d_n)^{d_n} \\ &\leq 2 \exp\left(-\log \frac{d_n}{\sqrt{n}} + \kappa\right) \rightarrow 0 \text{ as } n \rightarrow \infty. \end{aligned}$$

This shows that  $\text{TV}(U_n; [0, T]) \rightarrow 0$  almost surely.

Next, since  $d_n(1 - \lambda_n) \rightarrow \infty$ , Lemma 31 shows that

$$\boldsymbol{\mu}_n \rightarrow \boldsymbol{f}_1 \in l_1^\downarrow \text{ as } n \rightarrow \infty. \quad (7.101)$$

Therefore  $\boldsymbol{G}_n(0) = \boldsymbol{\mu}_n + \frac{\boldsymbol{Z}_n(0)}{\sqrt{n}} \rightarrow \boldsymbol{f}_1$  in  $l_1^\downarrow$ . Then by Lemma 27,

$$\sqrt{n}\boldsymbol{M}_n \Rightarrow \sqrt{2}B\boldsymbol{e}_1 \text{ in } \mathbb{D}([0, \infty) : l_2), \quad (7.102)$$

and by Lemma 29,  $\|\boldsymbol{Z}_{n,r+}\|_{2,T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . Define

$$V_n(t) \doteq - \int_0^t t_{n,2}(Z_{n,2}(s)) ds + \sqrt{n}M_{n,2}(t) - U_n(t).$$

Using (7.100) and Lemma 26 once more, we see that (7.97) is satisfied. Finally, for  $i \in \{3, \dots, r\}$ , define

$$W_{n,i}(t) \doteq \int_0^t t_{n,i-1}(Z_{n,i-1}(s)) ds - \int_0^t t_{n,i}(Z_{n,i}(s)) + \sqrt{n}M_{n,i}(t).$$

Then, from Lemma 26 again, it follows that (7.98) is satisfied with the above choice of  $W_{n,i}$ . Lemma 28 along with (7.101), (7.102), and  $|Z_{n,i}|_{*,T \wedge \tau_{n,L}} \leq L$  show that, as  $n \rightarrow \infty$ ,  $|V_n|_{*,T_n} \xrightarrow{P} 0$  and  $|W_{n,i}|_{*,T_n} \xrightarrow{P} 0$  for each  $i \in \{3, \dots, r\}$ .  $\square$

**Corollary 13.** *Suppose that the assumptions in Lemma 41 are satisfied. Assume further that  $d_n \ll n^{2/3}$ . Then, the conclusions of Lemma 41 hold with  $\theta_n = \alpha_n$  and*

$$\eta_n(t) \doteq \int_0^t \gamma_n^{-1} (1 + \delta_n(s))^+ e^{\gamma_n(Z_{n,1}(s) - \alpha_n)} \mathbb{I}_{\{Z_{n,1}(s) \geq \alpha_n\}} ds, \quad (7.103)$$

where  $\gamma_n \doteq \frac{d_n}{\sqrt{n}}$  and  $\delta_n$  is a process with sample paths in  $\mathbb{D}([0, \infty), \mathbb{R})$  such that  $|\delta_n|_{*, T \wedge \tau_{n,L}} \rightarrow 0$  a.s. for each  $L > 0$ .

*Proof.* Since  $d_n \ll n^{2/3}$  and  $\alpha_n = O(n^{1/6})$ , the hypothesis of Lemma 32 is satisfied. Define

$$\delta_n(s) \doteq t_{n,1}(Z_{n,1}(s)) \gamma_n \left( e^{\gamma_n[Z_{n,1}(s) - \alpha_n]} - e^{-\gamma_n \alpha_n} \right)^{-1} - 1.$$

Since  $\sup_{s \leq T \wedge \tau_{n,L}} |Z_{n,1}(s)| \leq L$ , Lemma 32 shows that  $|\delta_n|_{*, T_n} \rightarrow 0$  a.s. as  $n \rightarrow \infty$ . Next define

$$\begin{aligned} U_n(s) &\doteq \gamma_n^{-1} \int_0^s (1 + \delta_n(s)) \left( e^{-\gamma_n \alpha_n} - e^{\gamma_n(Z_{n,1}(s) - \alpha_n)} \mathbb{I}_{\{Z_{n,1}(s) < \alpha_n\}} \right) ds \\ &\quad + \int_0^s \gamma_n^{-1} (1 + \delta_n(s))^- e^{\gamma_n(Z_{n,1}(s) - \alpha_n)} \mathbb{I}_{\{Z_{n,1}(s) \geq \alpha_n\}} ds \end{aligned}$$

Then  $U_n(0) = 0$ ,  $U_n$  is absolutely continuous and, with  $\kappa = \sup_n \frac{d_n}{\sqrt{n}} \alpha_n^- < \infty$ ,

$$\begin{aligned} \text{TV}(U_n; [0, T_n]) \mathbb{I}_{\{|\delta_n|_{*, T_n} < 1\}} &= \gamma_n^{-1} \int_0^{T_n} |1 + \delta_n(s)| \left| e^{-\gamma_n \alpha_n} - e^{\gamma_n(Z_{n,1}(s) - \alpha_n)} \mathbb{I}_{\{Z_{n,1}(s) < \alpha_n\}} \right| ds \\ &\leq \frac{2(1 + e^\kappa)T}{\gamma_n} \rightarrow 0 \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Hence, since  $|\delta_n|_{*, T_n} \rightarrow 0$ , we have that  $\text{TV}(U_n; [0, T_n]) \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . By rearranging terms we see that, with the above definitions of  $U_n$  and  $\eta_n$ , (7.100) is satisfied. The result follows.  $\square$

Since  $\gamma_n \rightarrow \infty$  and  $\theta_n \rightarrow \alpha$  as  $n \rightarrow \infty$ , the previous lemma suggests a connection to the Skorokhod map  $\Gamma_\alpha$  defined in (7.2). In order to make this connection precise, we begin with the following lemma.

**Lemma 42.** *Under the assumptions of Theorem 11, for any  $L \in (0, \infty)$*

$$\sup_{t \in [0, T \wedge \tau_{n,L}]} (Z_{n,1}(t) - \alpha_n)^+ \xrightarrow{P} 0 \text{ as } n \rightarrow \infty. \quad (7.104)$$

*Proof.* Consider first the case  $d_n \gg \sqrt{n} \log n$ . For this case  $\epsilon_n \doteq \frac{\sqrt{n} \log d_n}{d_n} \rightarrow 0$ , and since

$$Z_{n,1}(t) = \sqrt{n}(G_{n,1}(t) - \lambda_n) \leq \sqrt{n}(1 - \lambda_n) = \frac{\sqrt{n} \log d_n}{d_n} + \alpha_n,$$

we have that (7.104) holds. Now consider the complementary case, namely  $d_n = O(\sqrt{n} \log n)$ . We will use Corollary 13. Since  $Z_{n,1}(0) \xrightarrow{P} z_1 \in \mathbb{R}$  with  $z_1 \leq \alpha$ , we have  $(Z_{n,1}(0) - \alpha_n)^+ \xrightarrow{P} 0$  as  $n \rightarrow \infty$ . It now suffices to show that for any  $\epsilon > 0$

$$\mathbf{P}\left(\sup_{t \in [0, T \wedge \tau_{n,L}]} Z_{n,1}(t) > \alpha_n + 6\epsilon\right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Let  $\tau_n \doteq \inf\{t \geq 0 \mid Z_{n,1}(t) > \alpha_n + 6\epsilon\}$  and, as before,  $T_n \doteq T \wedge \tau_{n,L}$ . It is then enough to show that  $\mathbf{P}(\tau_n \leq T_n) \rightarrow 0$  as  $n \rightarrow \infty$ . For this inductively define stopping times,  $\sigma_{n,0} = 0$ ,

$$\begin{aligned} \sigma_{n,2k-1} &= \inf\{t > \sigma_{n,2k-2} \mid Z_{n,1}(t) > \alpha_n + 3\epsilon\}, \quad k \in \mathbb{N}. \\ \sigma_{n,2k} &= \inf\{t > \sigma_{n,2k-1} \mid Z_{n,1}(t) < \alpha_n + 2\epsilon\} \end{aligned}$$

Note that for each  $n \in \mathbb{N}$ ,  $\sigma_{n,r} \rightarrow \infty$  as  $r \rightarrow \infty$ , almost surely. Also, henceforth, without loss of generality we consider only  $n$  that are large enough so that  $1/\sqrt{n} < \epsilon$ . Hence on the set  $\tau_n < \infty$ ,  $\tau_n \in [\sigma_{n,2k-1}, \sigma_{n,2k})$  for some  $k \in \mathbb{N}$ . Then for every  $K \in \mathbb{N}$

$$\mathbf{P}(\tau_n \leq T_n) \leq \sum_{k=1}^K \mathbf{P}(\tau_n \in [\sigma_{n,2k-1}, \sigma_{n,2k} \wedge T_n]) + \mathbf{P}(\sigma_{n,2K+1} \leq T_n).$$

Hence to complete the proof it is enough to show that

1. For each  $k \in \mathbb{N}$ ,  $\lim_{n \rightarrow \infty} \mathbf{P}(\tau_n \in [\sigma_{n,2k-1}, \sigma_{n,2k} \wedge T_n]) = 0$ ,
2.  $\lim_{K \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbf{P}(\sigma_{n,2K+1} \leq T_n) = 0$ .

Consider (1) first. Note that on the set  $C_{n,1} \doteq \{Z_{n,1}(0) \leq \alpha_n + 3\epsilon\}$ , for any  $k \in \mathbb{N}$ ,

$$\begin{aligned} \alpha_n + 3\epsilon &\leq Z_{n,1}(\sigma_{n,2k-1}) = Z_{n,1}(\sigma_{n,2k-1}-) + Z_{n,1}(\sigma_{n,2k-1}) - Z_{n,1}(\sigma_{n,2k-1}-) \\ &\leq Z_{n,1}(\sigma_{n,2k-1}-) + \epsilon \leq \alpha_n + 4\epsilon. \end{aligned} \tag{7.105}$$

Similarly,

$$Z_{n,1}(t) \geq \alpha_n + \epsilon \quad \text{for each } t \in [\sigma_{n,2k-1}, \sigma_{n,2k}]. \quad (7.106)$$

Let  $M'_n(t) \doteq \sqrt{n}M_{n,1}(t + \sigma_{n,2k-1}) - \sqrt{n}M_{n,1}(\sigma_{n,2k-1})$  for  $t \geq 0$  and consider the sets

$$C_{n,2} \doteq \{\tau_n, \sigma_{n,2k-1} \leq T_n\}, \quad C_{n,3} \doteq \left\{ |U_n|_{*,T_n} \leq \epsilon/2, |\delta_n|_{*,T_n} \leq \frac{1}{2} \right\}.$$

Then on the set  $C_n = \cap_{i=1}^3 C_{n,i}$ , using Corollary 13, for any  $t \in [0, (T_n \wedge \sigma_{n,2k}) - \sigma_{n,2k-1}]$ ,

$$\begin{aligned} & Z_{n,1}(t + \sigma_{n,2k-1}) - Z_{n,1}(\sigma_{n,2k-1}) \\ &= - \int_{\sigma_{n,2k-1}}^{\sigma_{n,2k-1}+t} (Z_{n,1}(s) - Z_{n,2}(s)) ds + M'_n(t) + U_n(t + \sigma_{n,2k-1}) - U_n(\sigma_{n,2k-1}) \\ & \quad - \int_{\sigma_{n,2k-1}}^{\sigma_{n,2k-1}+t} \gamma_n^{-1} (1 + \delta_n(s))^+ e^{\gamma_n(Z_{n,1}(s) - \alpha_n)} \mathbb{I}_{\{Z_{n,1}(s) \geq \alpha_n\}} ds \end{aligned}$$

Since for  $t$  in the above interval  $\sigma_{n,2k-1} + t \leq T_n \leq \tau_{n,L}$ ,  $|Z_{n,1}(s)| + |Z_{n,2}(s)| \leq 2L$  for any  $s \leq \sigma_{n,2k-1} + t$ . Also, since  $\sigma_{n,2k-1} + t \leq \sigma_{n,2k}$ , by (7.106),  $Z_{n,1}(s) - \alpha_n \geq \epsilon$  for any  $s \in [\sigma_{n,2k-1}, \sigma_{n,2k-1} + t]$ . Thus on  $C_n$  we have

$$Z_{n,1}(t + \sigma_{n,2k-1}) - Z_{n,1}(\sigma_{n,2k-1}) \leq 2Lt + M'_n(t) + \epsilon - \frac{t}{2\gamma_n} \exp(\gamma_n \epsilon) \doteq Y_n(t). \quad (7.107)$$

Using (7.105), on  $C_n$ ,  $Z_n(\tau_n) - Z_n(\sigma_{n,2k-1}) \geq \alpha_n + 6\epsilon - \alpha_n - 4\epsilon = 2\epsilon$ . Hence

$$\begin{aligned} \mathbf{P}(\tau_n \in [\sigma_{n,2k-1}, \sigma_{n,2k} \wedge T_n]) &\leq \mathbf{P}(\tau_n \in [\sigma_{n,2k-1}, \sigma_{n,2k} \wedge T_n], C_n) + \mathbf{P}(C_{n,1}^c) + \mathbf{P}(C_{n,3}^c) \\ &\leq \mathbf{P}\left(\sup_{t \in [0, T]} Y_n(t) \geq 2\epsilon\right) + \mathbf{P}(C_{n,1}^c) + \mathbf{P}(C_{n,3}^c), \end{aligned} \quad (7.108)$$

where the second inequality is on observing that on the set  $\{\tau_n \in [\sigma_{n,2k-1}, \sigma_{n,2k} \wedge T_n]\}$ , (7.107) holds with  $t$  replaced by  $\tau_n$ . Next note that  $M'_n$  is a  $\{\mathcal{G}_t^n\}$  martingale, where  $\mathcal{G}_t^n = \mathcal{F}_{t+\sigma_{n,2k-1}}^n$  and

$$\begin{aligned} \langle M'_n \rangle_t &= \langle \sqrt{n}M_{n,1} \rangle_{t+\sigma_{n,2k-1}} - \langle \sqrt{n}M_{n,1} \rangle_{\sigma_{n,2k-1}} \\ &= \int_{\sigma_{n,2k-1}}^{\sigma_{n,2k-1}+t} [G_{n,1}(s) - G_{n,2}(s) ds + \lambda_n(1 - \beta_n(G_{n,1}(s)))] ds \\ &\leq 2t, \end{aligned}$$

where the second equality is from (7.22).

Since  $\gamma_n \rightarrow \infty$  we can apply Lemma 34 to conclude

$$\mathbf{P}\left(\sup_{t \in [0, T]} Y_n(t) \geq 2\epsilon\right) = \mathbf{P}\left(\sup_{t \in [0, T]} \left[M'_n(t) - \left(\frac{\exp(\gamma_n \epsilon)}{2\gamma_n} - 2L\right)t\right] \geq \epsilon\right) \rightarrow 0$$

as  $n \rightarrow \infty$ . We also have  $\lim_n \mathbf{P}(C_{n,i}^c) = 0$  for  $i = 1, 3$  since, as noted earlier  $(Z_{n,1}(0) - \alpha_n)^+ \xrightarrow{P} 0$ , and by Corollary 13, respectively. From these observations it follows that the right side of (7.108) converges to 0 as  $n \rightarrow \infty$ , which completes the proof of (1)

Now we prove (2). Let  $\rho_i \doteq \sigma_{n,i} \wedge \tau_{n,L}$  and define

$$Y_{n,K}(t) \doteq \sum_{i=0}^K (Z_{n,1}(t \wedge \rho_{n,2i+1}) - Z_{n,1}(t \wedge \rho_{n,2i})).$$

Note that  $\{\sigma_{n,2K+1} \leq T_n\} \subseteq \{Y_{n,K}(T) \geq K\epsilon\}$  and hence to prove (2) it is sufficient to show that

$$\limsup_{n \rightarrow \infty} \mathbf{P}(Y_{n,K}(T) \geq K\epsilon) \rightarrow 0 \text{ as } K \rightarrow \infty. \quad (7.109)$$

From Corollary 13, we have that on the set  $C_{n,4} \doteq \{\text{TV}(U_n; [0, T_n]) \leq 1\}$ ,

$$\begin{aligned} Y_{n,K}(T) &= \sum_{i=0}^K \int_{T \wedge \rho_{n,2i}}^{T \wedge \rho_{n,2i+1}} (Z_{n,2}(s) - Z_{n,1}(s)) ds + \sum_{i=0}^K \sqrt{n} M_{n,1}(T \wedge \rho_{n,2i+1}) - \sqrt{n} M_{n,1}(T \wedge \rho_{n,2i}) \\ &\quad + \sum_{i=0}^K U_n(T \wedge \rho_{n,2i+1}) - U_n(T \wedge \rho_{n,2i}) - \sum_{i=0}^K \int_{T \wedge \rho_{n,2i}}^{T \wedge \rho_{n,2i+1}} \gamma_n^{-1} (1 + \delta_n(s))^+ e^{\gamma_n(Z_{n,1}(s) - \alpha_n)} \mathbb{I}_{\{Z_{n,1}(s) > \alpha_n\}} ds \\ &\leq 2LT + \sum_{i=0}^K (\sqrt{n} M_{n,1}(T \wedge \rho_{n,2i+1}) - \sqrt{n} M_{n,1}(T \wedge \rho_{n,2i})) + \text{TV}(U_n; [0, T]) \\ &\leq 2LT + 1 + M'_{n,K}(T) \end{aligned}$$

where we have used the facts that  $\sup_{s \leq \tau_{n,L}} |Z_{n,1}(s)| \leq L$ , and that the rightmost term in the third line is non-positive. Also, here

$$M'_{n,K}(t) \doteq \sum_{i=0}^K (\sqrt{n} M_{n,1}(t \wedge \rho_{n,2i+1}) - \sqrt{n} M_{n,1}(t \wedge \rho_{n,2i})).$$

Using (7.22), we see that  $M'_{n,K}$  is a  $\mathcal{F}_t^n$ -martingale with quadratic variation given by

$$\begin{aligned}\langle M'_{n,K} \rangle_t &= \sum_{i=0}^K \left( \langle \sqrt{n}M_{n,1} \rangle_{t \wedge \rho_{n,2i+1}} - \langle \sqrt{n}M_{n,1} \rangle_{t \wedge \rho_{n,2i}} \right) \\ &= \sum_{i=0}^K \int_{t \wedge \rho_{n,2i}}^{t \wedge \rho_{n,2i+1}} (G_{n,1}(s) - G_{n,2}(s) + \lambda_n - \lambda_n \beta_n(G_{n,1}(s))) ds \\ &\leq 2t.\end{aligned}$$

Hence

$$\begin{aligned}\mathbf{P}(Y_{n,K}(T) \geq K\epsilon) &\leq \mathbf{P}(Y_{n,K}(T) \geq K\epsilon, C_{n,4}) + \mathbf{P}(C_{n,4}^c) \\ &\leq \mathbf{P}(M'_{n,K}(T) > K\epsilon - (2LT + 1)) + \mathbf{P}(C_{n,4}^c) \\ &\leq \frac{\mathbf{E}M_{n,K}^{\prime 2}(T)}{(K\epsilon - (2LT + 1))^2} + \mathbf{P}(C_{n,4}^c) \\ &\leq \frac{2T}{(K\epsilon - (2LT + 1))^2} + \mathbf{P}(C_{n,4}^c).\end{aligned}$$

From Corollary 13,  $\mathbf{P}(C_{n,4}^c) \rightarrow 0$  as  $n \rightarrow \infty$ . This together with the above display shows  $\lim_{K \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbf{P}(Y_{n,K}(T) \geq K\epsilon) = 0$ . Thus we have shown (7.109) and the proof of (2) is complete. The result follows.  $\square$

**Lemma 43.** *Suppose the hypothesis of Theorem 11 holds, then for each  $n \in \mathbb{N}$ , there is a real constant  $\theta_n = \alpha_n + O(\sqrt{n}/d_n) \geq 0$  and processes  $U'_n, V_n$  with sample paths in  $\mathbb{D}([0, \infty) : \mathbb{R})$  such that with  $\tilde{Z}_{n,1} \doteq Z_{n,1} \wedge \theta_n$*

$$\begin{aligned}\tilde{Z}_{n,1}(t) &= \Gamma_{\theta_n} \left( \tilde{Z}_{n,1}(0) - \int_0^\cdot (\tilde{Z}_{n,1}(s) - Z_{n,2}(s)) + \sqrt{n}M_{n,1}(\cdot) + U'_n(\cdot) \right)(t), \text{ and} \\ Z_{n,2}(t) &= Z_{n,2}(0) - \int_0^t (Z_{n,2}(s) - Z_{n,3}(s))ds + V_n(t) + \eta_n(t) \quad \text{for all } t > 0,\end{aligned} \tag{7.110}$$

where

$$\eta_n = \hat{\Gamma}_{\theta_n} \left( \tilde{Z}_{n,1}(0) - \int_0^\cdot (\tilde{Z}_{n,1}(s) - Z_{n,2}(s)) + \sqrt{n}M_{n,1}(\cdot) + U'_n(\cdot) \right). \tag{7.111}$$

Furthermore, for any  $L, T \in (0, \infty)$ ,  $|U'_n|_{*, T \wedge \tau_{n,L}}, |V_n|_{*, T \wedge \tau_{n,L}}$  and  $|(Z_{n,1} - \theta_n)^+|_{*, T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ .

*Proof.* Let  $\theta_n$  be as in Lemma 41. Since  $d_n \gg \sqrt{n}$ ,  $\theta_n = \alpha_n + o(1)$  and Lemma 42 shows

$$|(Z_{n,1} - \theta_n)^+|_{*,T \wedge \tau_{n,L}} \rightarrow 0. \quad (7.112)$$

Note that  $\tilde{Z}_{n,1} = Z_{n,1} - (Z_{n,1} - \theta_n)^+$ . Hence we can rewrite (7.96) and (7.97) as

$$\begin{aligned} \tilde{Z}_{n,1}(t) &= \tilde{Z}_{n,1}(0) - \int_0^t \tilde{Z}_{n,1}(s) ds + \int_0^t Z_{n,2}(s) ds + \sqrt{n} M_{n,1}(t) + U'_n(t) - \eta_n(t) \\ Z_{n,2}(t) &= Z_{n,2}(0) - \int_0^t (Z_{n,2}(s) - Z_{n,3}(s)) ds + V_n(t) + \eta_n(t), \end{aligned} \quad (7.113)$$

where

$$U'_n(t) \doteq U_n(t) - \int_0^t (Z_{n,1}(s) - \theta_n)^+ ds - (Z_{n,1}(t) - \theta_n)^+ + (Z_{n,1}(0) - \theta_n)^+.$$

The properties of  $\eta_n$  from Lemma 41 (and Corollary 13) say that  $\eta_n$  is a non-decreasing process, with  $\eta_n(0) = 0$  and  $\eta_n(t) = \int_0^t \mathbb{I}_{\{\tilde{Z}_{n,1}(s) = \theta_n\}} d\eta_n(s)$ . Since  $\tilde{Z}_{n,1} \leq \theta_n$ , (7.113) and the characterizing properties of the Skorokhod map show (7.110) and (7.111). Finally, by Lemma 41, Corollary 13, and Lemma 42

$$|U_n|_{*,T \wedge \tau_{n,L}} \xrightarrow{P} 0, \text{ and } |V_n|_{*,T \wedge \tau_{n,L}} \xrightarrow{P} 0$$

as  $n \rightarrow \infty$ . Hence, using (7.112),  $|U'_n|_{*,T \wedge \tau_{n,L}} \xrightarrow{P} 0$  as  $n \rightarrow \infty$ , and the result follows.  $\square$

The following lemma will be needed in order to prove the tightness of  $\mathbf{Z}_n$ .

**Lemma 44.** *Under the hypothesis of Theorem 11, the collection of random variables  $\left\{ \|\mathbf{Z}_n\|_{2,T} \right\}_{n \in \mathbb{N}}$  is tight for any  $T \in (0, \infty)$ .*

*Proof.* Fix  $T \in (0, \infty)$ . In Lemma 43 using the definition of the Skorokhod map  $\Gamma_{\theta_n}$  for  $\theta_n \geq 0$  (see (7.2)), we have, for any  $t > 0$  that

$$\eta_n(t) \leq \left| \tilde{Z}_{n,1}(0) \right| + \int_0^t \left| \tilde{Z}_{n,1}(s) \right| ds + \int_0^t |Z_{n,2}(s)| ds + |\sqrt{n} M_{n,1}|_{*,t} + |U'_n|_{*,t}.$$



This shows that for any  $t \geq 0$

$$\begin{aligned} |\tilde{Z}_{n,1}|_{*,t} &\leq 2 \left( |\tilde{Z}_{n,1}(0)| + \int_0^t |\tilde{Z}_{n,1}|_{*,s} ds + \int_0^t |Z_{n,2}|_{*,s} ds + |\sqrt{n}M_{n,1}|_{*,t} + |U'_n|_{*,t} \right) \\ |Z_{n,2}|_{*,t} &\leq |\tilde{Z}_{n,1}(0)| + |Z_{n,2}(0)| + \int_0^t |\tilde{Z}_{n,1}|_{*,s} ds + \int_0^t (2|Z_{n,2}|_{*,s} + |Z_{n,3}|_{*,s}) ds \\ &\quad + |\sqrt{n}M_{n,1}|_{*,t} + |U'_n|_{*,t} + |V_n|_{*,t}, \end{aligned}$$

and

$$|Z_{n,i}|_{*,t} \leq |Z_{n,i}(0)| + \int_0^t |Z_{n,i}|_{*,s} ds + \int_0^t |Z_{n,i+1}|_{*,s} ds + |W_{n,i}|_{*,t} \quad \text{for } i \in \{3, \dots, r\}$$

where the last line is from Lemma 41. Let  $H_t \doteq |\tilde{Z}_{n,1}|_{*,t} + |Z_{n,2}|_{*,t} + \dots + |Z_{n,r}|_{*,t}$ . By adding over equations in the above display, we have for  $t \in [0, \tau]$  and  $\tau \in [0, T]$  that

$$0 \leq H_t \leq 4 \left( H_0 + |\sqrt{n}M_{n,1}|_{*,\tau} + |U'_n|_{*,\tau} + |V_n|_{*,\tau} + \sum_{i=3}^r |W_{n,i}|_{*,\tau} + \int_0^t H_s ds \right).$$

By Gronwall's inequality, for all  $\tau \in [0, T]$ ,

$$H_\tau \leq 4 \left( H_0 + |\sqrt{n}M_{n,1}|_{*,\tau} + |U'_n|_{*,\tau} + |V_n|_{*,\tau} + \sum_{i=3}^r |W_{n,i}|_{*,\tau} \right) e^{4\tau}. \quad (7.114)$$

Let  $\vec{Z}_n \doteq (\tilde{Z}_{n,1}, Z_{n,2}, \dots, Z_{n,r})$ . Since  $\vec{Z}_n(0) \xrightarrow{P} (z_1, \dots, z_r)$ , and  $\sqrt{n}\mathbf{M}_n \Rightarrow B\mathbf{e}_1$ , for every  $\epsilon > 0$  there is a  $L_1 \in (0, \infty)$  such that for every  $n \in \mathbb{N}$

$$P(C_{n,1}) \leq \frac{\epsilon}{2}, \text{ where } C_{n,1} \doteq \left\{ H_0 + |\sqrt{n}M_{n,1}|_{*,T} \geq L_1 \right\}.$$

Applying Lemmas 41 and Lemma 43 with  $L = 4(L_1 + 1)e^{4T} + 2$  we can find an  $n_0 \in \mathbb{N}$  so that  $P(C_{n,2}) \leq \frac{\epsilon}{2}$  for  $n \geq n_0$ , where

$$C_{n,2} \doteq \left\{ |U'_n|_{*,T_n} + |V_n|_{*,T_n} + \sum_{i=3}^r |W_{n,i}|_{*,T_n} + |(Z_{n,1} - \alpha_n)^+|_{*,T_n} + \|\mathbf{Z}_{n,r+}\|_{2,T_n} \geq 1 \right\}$$

and  $T_n \doteq T \wedge \tau_{n,L}$ . On the event  $(C_{n,1} \cup C_{n,2})^c$

$$\|\vec{Z}_n\|_{1,T_n} = H_{T_n} < 4(L_1 + 1)e^{4T}$$

by (7.114), and hence by triangle inequality

$$\begin{aligned}\|\mathbf{Z}_n\|_{2,T_n} &\leq \left\|\tilde{\mathbf{Z}}_n\right\|_{1,T_n} + |(Z_{n,1} - \alpha)^+|_{*,T_n} + \|\mathbf{Z}_{n,r+}\|_{2,T_n} \\ &< 4(L_1 + 1)e^{4T} + 1 = L - 1.\end{aligned}\tag{7.115}$$

Also, by the definition of  $\tau_{n,L}$ ,  $\|\mathbf{Z}_n(\tau_{n,L})\|_2 \geq L - \frac{1}{\sqrt{n}}$  on the set  $\tau_{n,L} < T$ . Hence we must have that  $\tau_{n,L} > T$  whenever (7.115) holds, and hence

$$\|\mathbf{Z}_n\|_{2,T} < L - 1 \text{ on the event } (C_{n,1} \cup C_{n,2})^c.$$

This shows that

$$\mathbf{P}\left(\|\mathbf{Z}_n\|_{2,T} \geq L\right) \leq \mathbf{P}(C_{n,1} \cup C_{n,2}) \leq \epsilon \quad \forall n \geq n_0$$

Since  $\epsilon > 0$  is arbitrary, the result follows.  $\square$

The following result is immediate from Lemmas 30, 41, 43, and 44.

**Corollary 14.** *Under the hypothesis of Theorem 11, for any  $T > 0$ ,  $\lim_{L \rightarrow \infty} \sup_n \mathbf{P}(\tau_{n,L} \leq T) = 0$ . In particular the processes  $W_{n,i}, U'_n, V_n, \|\mathbf{Z}_{n,r+}\|_2, (Z_{n,1} - \alpha_n)^+, (Z_{n,1} - \theta_n)^+$  converge in probability to zero in  $\mathbb{D}([0, \infty) : \mathbb{R})$ .*

**Corollary 15.** *Under the hypothesis of Theorem 11, the sequence of processes  $\{\mathbf{Z}_n\}_{n \in \mathbb{N}}$  is tight in  $\mathbb{D}([0, \infty) : l_2)$ .*

*Proof.* Let  $\theta_n$  be as in Lemma 43. Then by Corollary 14, for each fixed  $T < \infty$ ,  $\|\mathbf{Z}_{n,r+}\|_{2,T} \xrightarrow{P} 0$  and  $|(Z_{n,1} - \theta_n)^+|_{*,T} \xrightarrow{P} 0$ . Hence it is sufficient to show that the sequence  $\{\tilde{\mathbf{Z}}_n\}_{n \in \mathbb{N}}$  introduced in the proof of Lemma 44 is tight in  $\mathbb{D}([0, T] : \mathbb{R}^r)$ . From Lemma 44, the convergence of  $W_{n,i}$  in Corollary 14, and equations for  $Z_{n,j}$ ,  $j = 3, \dots, r$  in Lemma 41, it is immediate that  $(Z_{n,3}, \dots, Z_{n,r})$  is tight in  $\mathbb{D}([0, \infty) : \mathbb{R}^{r-2})$ . Finally consider the pair  $(\tilde{Z}_{n,1}, Z_{n,2})$ . Once again using Lemma 44, the convergence of  $\sqrt{n}M_{n,1}$  in Lemma 41, and the convergence of  $U'_n$  in Corollary 14, it follows that

$$R_n \doteq \tilde{Z}_{n,1}(0) - \int_0^\cdot \left( \tilde{Z}_{n,1}(s) - Z_{n,2}(s) \right) + \sqrt{n}M_{n,1}(\cdot) + U'_n(\cdot) \tag{7.116}$$

is tight in  $\mathbb{D}([0, \infty) : \mathbb{R})$ . Using the identity

$$\Gamma_{\theta_n}(R_n)(t) = \Gamma_{\theta_n}(\Gamma_{\theta_n}(R_n)(s) + R_n(\cdot + s) - R_n(s))(t - s)$$

for  $0 \leq s \leq t \leq T$ , we see from the definition of the Skorohod map that

$$|\Gamma_{\theta_n}(R_n)(t) - \Gamma_{\theta_n}(R_n)(s)| \leq 2 \sup_{s \leq u \leq t} |R_n(u) - R_n(s)|.$$

Together with the tightness of  $R_n$  this immediately implies the tightness of  $\tilde{Z}_{n,1} = \Gamma_{\theta_n}(R_n)$  and of  $\hat{\Gamma}_{\theta_n}(R_n)$ . Finally the tightness of  $Z_{n,2}$  is now immediate from Lemma 44, the convergence of  $V_n$  in Corollary 14 and the tightness of  $\hat{\Gamma}_{\theta_n}(R_n)$  noted above. The result follows.  $\square$

*Proof of Theorem 11.* From Lemma 31 and from the tightness of  $\{\|\mathbf{Z}_n(0)\|_1\}_{n \in \mathbb{N}}$ , it follows under the conditions of the theorem that  $\boldsymbol{\mu}_n \xrightarrow{P} \mathbf{f}_1$  and  $\mathbf{G}_n(0) \xrightarrow{P} \mathbf{f}_1$  in  $l_1^\downarrow$ . This proves the first statement in the theorem. Now consider the second statement. Fix  $T < \infty$ . From Corollary 15,  $\{\mathbf{Z}_n\}_{n \in \mathbb{N}}$  is tight in  $\mathbb{D}([0, \infty) : l_2)$ . Also from Lemma 41,  $\sqrt{n}M_{n,1}$  converges in distribution to  $\sqrt{2}B$  where  $B$  is a standard Brownian motion and from Corollary 14

$$(\{W_{n,i}\}_{i=3}^r, U'_n, V_n, (Z_{n,1} - \theta_n)^+) \xrightarrow{P} \mathbf{0} \text{ in } \mathbb{D}([0, T] : \mathbb{R}^r)$$

Suppose that along a subsequence

$$(\mathbf{Z}_n, \sqrt{n}M_{n,1}, \{W_{n,i}\}_{i=3}^r, U'_n, V_n, (Z_{n,1} - \theta_n)^+) \Rightarrow (\mathbf{Z}, \sqrt{2}B, \mathbf{0})$$

in  $\mathbb{D}([0, \infty) : l_2 \times \mathbb{R}^{r+2})$  and for notational simplicity label the subsequence once more as  $\{n\}$ . Also by appealing to Skorohod embedding theorem we assume that all the processes in the above display are given on a common probability space and the above convergence holds a.s. Since  $J(\mathbf{Z}_n) \leq \frac{1}{\sqrt{n}}$  and  $\mathbf{Z}_n(0) \xrightarrow{P} \mathbf{z}$ , we have  $J(\mathbf{Z}) = 0$  and  $\mathbf{Z}(0) = \mathbf{z}$  almost surely. In particular  $\mathbf{Z}$  has sample paths in  $C([0, \infty) : l_2)$  and  $(\mathbf{Z}_n, \sqrt{n}M_{n,1}) \rightarrow (\mathbf{Z}, \sqrt{2}B)$  uniformly over compact time intervals in  $l_2 \times \mathbb{R}$ . Since by Corollary 14, for every  $T < \infty$ ,  $\|\mathbf{Z}_{n,r+}\|_{2,T} \xrightarrow{P} 0$ , it suffices to show that  $(Z_1, \dots, Z_r)$  along with  $B$  satisfy (7.14).

From the equations of  $(Z_{n,3}, \dots, Z_{n,r})$  in Lemma 41, uniform convergence of  $Z_n$  to  $Z$ , and the uniform convergence of  $W_{n,i}$  to 0, it is immediate that  $(Z_3, \dots, Z_r)$  satisfy (7.14). Finally consider the equations for  $(Z_1, Z_2)$ . From (7.116) and uniform convergence properties observed above it is immediate that  $R_n$  converges uniformly, a.s., to  $R$  given as

$$R(\cdot) = Z_1(0) - \int_0^\cdot (Z_1(s) - Z_2(s)) + \sqrt{2}B(\cdot).$$

Since  $\theta_n = \alpha_n + O(\sqrt{n}/d_n) \rightarrow \alpha$ , this shows that, for every  $T < \infty$ ,

$$\begin{aligned} \Gamma_{\theta_n}(R_n)(t) &= R_n(t) - \sup_{s \in [0, t]} (R_n(t) - \theta_n)^+ \\ &\rightarrow R(t) - \sup_{s \in [0, t]} (R(t) - \alpha)^+ = \Gamma_\alpha(R)(t) \end{aligned}$$

uniformly for  $t \in [0, T]$ , a.s., where  $(R(t) - \alpha)^+$  is taken to be 0 when  $\alpha = \infty$ . Similarly,

$$\hat{\Gamma}_{\theta_n}(R_n)(t) \rightarrow \hat{\Gamma}_\alpha(R)(t)$$

uniformly for  $t \in [0, T]$ , a.s. Here, when  $\alpha = \infty$ ,  $\Gamma_\alpha$  and  $\hat{\Gamma}_\alpha$  are as introduced in (7.12). The fact that  $(Z_1, Z_2)$  solve the first two equations in (7.14) is now immediate from Lemma 43, the convergence  $\tilde{Z}_{n,1} - Z_{n,1} \xrightarrow{P} 0$ , and the uniform convergence of  $V_n$  to 0 noted previously. The result follows.  $\square$

## 7.10 Technical estimates

In this section we provide proofs of the various results in Section 7.5 and of Lemma 33.

### 7.10.1 Proof of Lemma 21:

*Proof.* Fix  $\epsilon \in (0, 1)$ . First suppose  $\frac{d_n}{n} \rightarrow 0$ . Consider  $x \in (\epsilon, 1]$ . Let  $\Delta_n(x) \doteq \log \beta_n(x) - \log \gamma_n(x)$ .

Let  $n_0 \in \mathbb{N}$  be such that for all  $n \geq n_0$ ,  $d_n/n < \epsilon/2$ . Then, for  $n \geq n_0$ ,

$$\begin{aligned} \Delta_n(x) &\doteq \sum_{i=0}^{d_n-1} \log \left( \frac{x - i/n}{1 - i/n} \right) - \log x^{d_n} = \sum_{i=0}^{d_n-1} \left\{ \log \left( \frac{x - i/n}{1 - i/n} \right) - \log x \right\}, \\ &= \sum_{i=0}^{d_n-1} \log \left( \frac{1 - i/(nx)}{1 - i/n} \right) = \sum_{i=0}^{d_n-1} \log \left( 1 - (i/n) \frac{1/x - 1}{1 - i/n} \right). \end{aligned} \quad (7.117)$$

Differentiating  $\Delta_n$  gives,

$$\Delta'_n(x) = \sum_{i=0}^{d_n-1} \left( \frac{1}{x - i/n} - \frac{1}{x} \right) = \sum_{i=0}^{d_n-1} \frac{i/n}{x(x - i/n)}.$$

Since  $n \geq n_0$  and  $x \in [\epsilon, 1]$  we have  $x(x - \frac{i}{n}) \geq \epsilon^2/2$  for  $i \leq d_n - 1$ . Hence,

$$|\Delta'_n(x)| \leq \frac{2}{\epsilon^2} \sum_{i=0}^{d_n-1} (i/n) \leq \frac{1}{\epsilon^2} \frac{d_n^2}{n}.$$

From the definition of  $\Delta_n$ , we also have,

$$\Delta'_n(x) = \frac{\beta'_n(x)}{\beta_n(x)} - \frac{\gamma'_n(x)}{\gamma_n(x)} = \frac{\gamma'_n(x)}{\gamma_n(x)} \left( \frac{\beta'_n(x)}{\gamma'_n(x)} \frac{\gamma_n(x)}{\beta_n(x)} - 1 \right). \quad (7.118)$$

Since  $\frac{\gamma'_n(x)}{\gamma_n(x)} = \frac{d_n}{x} \geq d_n$  for  $x \in [\epsilon, 1]$ , from (7.118) we have,

$$\sup_{x \in [\epsilon, 1]} \left| \frac{\beta'_n(x)}{\gamma'_n(x)} \frac{\gamma_n(x)}{\beta_n(x)} - 1 \right| \leq \frac{1}{d_n} \sup_{x \in [\epsilon, 1]} |\Delta'_n(x)| \leq \frac{1}{\epsilon} \frac{d_n}{n} \rightarrow 0.$$

This proves (7.39).

Now assume  $\frac{d_n}{\sqrt{n}} \rightarrow 0$ . Once more consider  $x \in (\epsilon, 1]$  and  $n \geq n_0$ . Let  $C \doteq \sup_{n \geq n_0} \frac{1/\epsilon - 1}{1 - d_n/n} < \infty$  and let  $n_1 > n_0$  be such that  $d_n C/n < 1/2$  for all  $n \geq n_1$ . Then for  $n \geq n_1$  and  $x \in [\epsilon, 1]$ :

$$|\Delta_n(x)| \leq \sum_{i=0}^{d_n-1} 2 \left| (i/n) \frac{1/x - 1}{1 - i/n} \right| \leq 2C \sum_{i=0}^{d_n-1} i/n \leq C \frac{d_n^2}{n}, \quad (7.119)$$

where the first inequality is from (7.117) and the inequality  $|\log(1+h)| \leq 2|h|$  for  $|h| \leq 1/2$ . This shows  $\sup_{x \in [\epsilon, 1]} |\Delta_n(x)| \rightarrow 0$ , hence showing the first convergence in (7.40). Finally the second convergence (7.40) is immediate on combining the first convergence with (7.39).  $\square$

### 7.10.2 Proof of Corollary 7:

This is an immediate consequence of the estimate in (7.119).

### 7.10.3 Proof of Corollary 8:

*Proof.* Let  $\epsilon > 0$  and  $n_0 \in \mathbb{N}$  be such that  $\mu_{n,i} > \epsilon$  for all  $n \geq n_0$ . By Lemma 21, as  $n \rightarrow \infty$

$$\frac{\beta'_n(\mu_{n,i})}{\beta_n(\mu_{n,i})} = (1 + o(1)) \frac{\gamma'_n(\mu_{n,i})}{\gamma_n(\mu_{n,i})}. \quad (7.120)$$

Recall that  $\mu_{n,i+1} \doteq \lambda_n \beta_n(\mu_{n,i})$  and  $\gamma_n(x) \doteq x^{d_n}$ . Hence (7.120) gives

$$\frac{\beta'_n(\mu_{n,i})}{\mu_{n,i+1}/\lambda_n} = (1 + o(1)) \frac{d_n}{\mu_{n,i}} \quad (7.121)$$

completing the proof.  $\square$

### 7.10.4 Proof of Lemma 22:

*Proof.* From Corollary 7, there is a  $n_0 \in \mathbb{N}$  and  $C \in (0, \infty)$  such that for all  $n \geq n_0$

$$\sup_{x \in [\epsilon, 1]} |\log \beta_n(x) - \log \gamma_n(x)| \leq \frac{C d_n^2}{n}.$$

Thus, if for  $n \geq n_0$  and  $i \in \mathbb{N}$ ,  $\mu_{n,i} \geq \epsilon$ , then

$$\begin{aligned} \log \mu_{n,i+1} &= \log \lambda_n + \log \beta_n(\mu_{n,i}) = \log \lambda_n + \log \gamma_n(\mu_{n,i}) + \gamma_{n,i} \\ &= \log \lambda_n + d_n \log \mu_{n,i} + \gamma_{n,i}, \end{aligned} \quad (7.122)$$

where  $|\gamma_{n,i}| \leq \frac{Cd_n^2}{n}$ . Now let  $k \in \mathbb{N}$  and  $n_1 \in \mathbb{N}$  be such that for all  $n \geq n_1$ ,  $\mu_{n,k} \geq \epsilon$ . We will show that for  $n \geq n_0 \vee n_1$  and  $j \in \{1, \dots, k\}$  that

$$\log \mu_{n,j+1} = (\log \lambda_n) \left( \sum_{i=0}^j d_n^i \right) + \beta_{n,j}, \quad (7.123)$$

where  $|\beta_{n,j}| \leq \frac{C}{n} \sum_{i=1}^j d_n^{i+1}$ . Note the the lemma is immediate from (7.123) on taking  $j = k$ . To prove (7.123) we argue inductively. First note that since  $\mu_n \in l_1^\perp$ ,  $\mu_{n,i} \geq \mu_{n,k} \geq \epsilon$  for each  $i \leq k$  and  $n \geq n_1$ . Hence (7.122) holds for each  $i \leq k$  and  $n \geq n_0 \vee n_1$ . Taking  $i = 1$  in (7.122) and noting that  $\mu_{n,1} = \lambda_n$  proves (7.123) for the case  $j = 1$ .

Suppose now (7.123) holds for some  $j \leq k - 1$ . Then, using  $i = j + 1$ , in (7.122)

$$\log \mu_{n,j+2} = \log \lambda_n + d_n \log \mu_{n,j+1} + \gamma_{n,j+1},$$

where  $|\gamma_{n,j+1}| \leq \frac{Cd_n^2}{n}$ . By the induction hypothesis, (7.123) holds for  $j$ . Hence

$$\log \mu_{n,j+2} = \log \lambda_n + d_n \left\{ (\log \lambda_n) \left( \sum_{i=0}^j d_n^i \right) + \beta_{n,j} \right\} = (\log \lambda_n) \left( \sum_{i=0}^{j+1} d_n^i \right) + d_n \beta_{n,j} + \gamma_{n,j}$$

and hence  $\beta_{n,j+1} = d_n \beta_{n,j} + \gamma_{n,j}$ . This shows

$$|\beta_{n,j+1}| = |d_n \beta_{n,j} + \gamma_{n,j}| \leq d_n \frac{C}{n} \sum_{i=1}^j d_n^{i+1} + \frac{Cd_n^2}{n} = \frac{C}{n} \sum_{i=1}^{j+1} d_n^{i+1}$$

which shows that (7.123) holds for  $j + 1$ . This completes the proof.  $\square$

### 7.10.5 Proof of Corollary 9:

*Proof.* Since  $d_n \rightarrow \infty$ , the assumption  $\frac{\xi_n^2}{d_n} \rightarrow 0$  shows that  $\frac{|\xi_n|}{d_n} \leq \frac{1+\xi_n^2}{d_n} \rightarrow 0$ . This shows that  $\epsilon_n \doteq 1 - \lambda_n = \frac{\xi_n + \log d_n}{d_n^k}$  also converges to 0.

We first show that  $\mu_{n,i} \rightarrow 1$  for each  $i \in \{1, \dots, k\}$ . We will argue inductively. Since  $\mu_{n,1} \doteq \lambda_n = 1 - \epsilon_n$ , we have  $\mu_{n,1} \rightarrow 1$ . Suppose now that  $\mu_{n,i} \rightarrow 1$  for some  $i \leq k - 1$ . Hence eventually  $\mu_{n,i} \geq \frac{1}{2}$ . Applying Lemma 22 with  $k = i$  and  $\epsilon = \frac{1}{2}$  and simplifying the resulting expression, we

get

$$\log \mu_{n,i+1} = (\log \lambda_n) \frac{d_n^{i+1} - 1}{d_n - 1} + O\left(\frac{d_n^2(d_n^i - 1)}{n(d_n - 1)}\right) \quad (7.124)$$

$$= O(\epsilon_n) \frac{d_n^{i+1} - 1}{d_n - 1} + O\left(\frac{d_n^2(d_n^i - 1)}{n(d_n - 1)}\right) = O\left(\frac{\xi_n + \log d_n}{d_n^{k-i}}\right) + O\left(\frac{d_n^{i+1}}{n}\right), \quad (7.125)$$

where the second equality uses  $\log \lambda_n = \log(1 - \epsilon_n) = O(\epsilon_n)$  and the third follows on recalling that  $d_n \rightarrow \infty$ . Since  $i \leq k - 1$ ,  $\frac{|\xi_n|}{d_n^{k-i}} \leq \frac{1 + \xi_n^2}{d_n} \rightarrow 0$ . Using this along with  $d_n^{k+1} \ll n$  in (7.125) shows that  $\mu_{n,i+1} \rightarrow 1$ . Hence, by induction,  $\mu_{n,i} \rightarrow 1$  for  $i \leq k$ .

Next we argue that  $\beta'_n(\mu_{n,k}) \rightarrow \alpha$ . Since  $\lambda_n \rightarrow 1$  and  $\mu_{n,k} \rightarrow 1$ , from Corollary 8 we have that

$$\lim_{n \rightarrow \infty} \frac{\beta'_n(\mu_{n,k})}{d_n \mu_{n,k+1}} = 1$$

Hence it suffices to show that  $d_n \mu_{n,k+1} \rightarrow \alpha$ . For this note that

$$\begin{aligned} \log(d_n \mu_{n,k+1}) &= \log \mu_{n,k+1} + \log d_n \\ &= \log(1 - \epsilon_n) \left( \frac{d_n^{k+1} - 1}{d_n - 1} \right) + O\left( \frac{d_n^2 d_n^k - 1}{n d_n - 1} \right) + \log d_n \\ &= (-\epsilon_n + O(\epsilon_n^2)) d_n^k (1 + O(1/d_n)) + \log d_n + O\left( \frac{d_n^{k+1}}{n} \right), \end{aligned}$$

where the second equality is from (7.124) and last equality is by using Taylor's expansion for  $\log(1 - \epsilon_n)$ . Using  $d_n^{k+1} \ll n$  and  $|\epsilon_n^2 d_n^k| \leq \frac{2(\xi_n^2 + (\log d_n)^2)}{d_n^k} \rightarrow 0$ , we now have

$$\begin{aligned} \log(d_n \mu_{n,k+1}) &= (-\epsilon_n d_n^k + o(1))(1 + O(1/d_n)) + \log d_n + o(1) \\ &= (-\xi_n - \log d_n)(1 + O(1/d_n)) + \log d_n + o(1) \\ &= -\xi_n - \log d_n + \log d_n + O\left( \frac{\xi_n + \log d_n}{d_n} \right) + o(1) = -\xi_n + o(1) \rightarrow \log(\alpha) \end{aligned}$$

where the last equality once more uses the observation that  $\frac{|\xi_n|}{d_n} \rightarrow 0$ . Thus we have  $d_n \mu_{n,k+1} \rightarrow \alpha$  as  $n \rightarrow \infty$  which completes the proof.  $\square$



### 7.10.6 Proof of Lemma 23:

*Proof.* Since  $\mu_{n,k} \rightarrow 1$  and  $j \mapsto \mu_{n,j}$  is nonincreasing, we have  $\mu_{n,i} \rightarrow 1$  for each  $i \leq k$ . Additionally, since  $\lambda_n \rightarrow 1$ , Corollary 8 shows that for any  $i \in [k]$   $\lim_{n \rightarrow \infty} \frac{\beta'_n(\mu_{n,i})}{d_n \mu_{n,i+1}} = 1$ . As a consequence,  $\beta'_n(\mu_{n,k-1}) \rightarrow \infty$  as  $n \rightarrow \infty$ , and for any  $j \in [k-2]$

$$\lim_{n \rightarrow \infty} \frac{\beta'_n(\mu_{n,j})}{\beta'_n(\mu_{n,j+1})} = \lim_{n \rightarrow \infty} \frac{d_n \mu_{n,j+1}}{d_n \mu_{n,j+2}} = \lim_{n \rightarrow \infty} \frac{\mu_{n,j+1}}{\mu_{n,j+2}} = 1.$$

This completes the proof of the lemma.  $\square$

### 7.10.7 Proof of Lemma 24:

*Proof.* By the first part of Lemma 21, (7.42) is immediate from (7.41). Now consider (7.41). Taking logarithms in (7.8), for  $x > d_n/n$ ,

$$\log \beta_n(x) = \sum_{i=0}^{d_n-1} \left( \log \left( x - \frac{i}{n} \right) - \log \left( 1 - \frac{i}{n} \right) \right) = \sum_{i=0}^{d_n-1} \left( \log \left( 1 - \frac{i}{n} - (1-x) \right) - \log \left( 1 - \frac{i}{n} \right) \right).$$

Let  $\delta_n = \epsilon_n + \frac{d_n}{n}$ . For large  $n$ ,  $\delta_n \leq \frac{1}{2}$ , and hence, using the expansion  $\log(1-h) = -h + O(h^2)$  for  $|h| \leq \frac{1}{2}$ , for any  $x \in [1 - \epsilon_n, 1]$ :

$$\begin{aligned} \log \beta_n(x) &= \sum_{i=0}^{d_n-1} \left\{ -\frac{i}{n} - (1-x) + \frac{i}{n} + O(\delta_n^2) \right\} = -d_n(1-x) + O(d_n \delta_n^2) \\ &= d_n \log(1 - (1-x)) + O(d_n \delta_n^2) = \log \gamma_n(x) + O(d_n \delta_n^2). \end{aligned}$$

Note that  $\delta_n^2 = (\epsilon_n + d_n/n)^2 \leq 2\left(\epsilon_n^2 + \frac{d_n^2}{n^2}\right)$ . Hence by our assumptions  $d_n \delta_n^2 \rightarrow 0$ . This proves (7.41) and completes the proof of the lemma.  $\square$

### 7.10.8 Proof of Lemma 25:

*Proof.* By (7.37)

$$\sup_{x \in [0, 1-\epsilon_n]} |\beta_n(x)| \leq (1 - \epsilon_n)^{d_n} = e^{-d_n \epsilon_n + o(1)} \rightarrow 0.$$

Similarly, by (7.38), under the assumption  $\limsup_n \frac{d_n}{n} < 1$ , for large  $n$ ,

$$\sup_{x \in [0, 1-\epsilon_n]} |\beta'_n(x)| \leq (1 - d_n/n)^{-1} d_n (1 - \epsilon_n)^{d_n-1} = e^{-d_n \epsilon_n + \log d_n + O(1)} \rightarrow 0.$$

□

### 7.10.9 Proof of Lemma 33

For a right continuous bounded variation function  $F : [0, T] \rightarrow \mathbb{R}$ , let  $dF$  denote the signed measure on  $(0, T]$  given by  $dF(a, b] = F(b) - F(a)$  for  $0 \leq a < b \leq T$ , and  $d\lambda$  denote the Lebesgue measure on  $(0, T]$ . Bounded measurable functions  $h : [0, T] \rightarrow \mathbb{R}$  act on signed measure  $d\mu$  on  $(0, T]$  on the left as follows:  $hd\mu$  denotes the signed measure  $A \mapsto \int_A h(x) d\mu(x)$ ,  $A \in \mathcal{B}(0, T]$ .

Let  $F(t) \doteq \int_0^t f(s) d\lambda(s)$  for  $t \in [0, T]$ . Note that  $z$  defined in (7.56) is a right continuous function with bounded variations. The corresponding measure  $dz$  on  $(0, T]$  satisfies the identity

$$dz = -fz d\lambda + gd\lambda + dM,$$

namely

$$dz + fz d\lambda = gd\lambda + dM.$$

Acting on the left in the above identity by the bounded continuous function  $e^F(x) \doteq e^{F(x)}$  we get

$$e^F dz + e^F fz d\lambda = e^F gd\lambda + e^F dM.$$

Since  $dF = f d\lambda$ , by the change of variable formula (cf. (86, Theorem VI.8.3))  $de^F = fe^F d\lambda$ . Hence

$$e^F dz + z de^F = e^F gd\lambda + e^F dM.$$

Two applications of the integration by parts formula (cf. (16, Theorem 18.4)) show that

$$d(e^F z) = e^F gd\lambda + d(e^F M) - M de^F.$$

Computing the total measure on  $(0, t]$  for  $t \leq T$ :

$$e^{F(t)}z(t) - z(0) = \int_0^t e^{F(s)}g(s)d\lambda(s) + e^{F(t)}M(t) - M(0) - \int_0^t M(s)de^F(s)$$

Rearranging terms and multiplying by  $e^{-F(t)}$  on both sides:

$$z(t) = \int_0^t e^{F(s)-F(t)}g(s)d\lambda(s) + M(t) - e^{-F(t)} \int_0^t M(s)de^F(s) + e^{-F(t)}(z(0) - M(0)). \quad (7.126)$$

We now estimate the various terms on the right hand side of (7.126). The first term on the right hand side of (7.126) satisfies for  $t \in [0, T \wedge \tau]$

$$\begin{aligned} \left| \int_0^t e^{F(s)-F(t)}g(s)d\lambda(s) \right| &\leq |g|_{*,T \wedge \tau} \int_0^t e^{F(s)-F(t)}d\lambda(s) \\ &\leq |g|_{*,T \wedge \tau} \int_0^t e^{-\int_s^t f(u)du}d\lambda(s) \\ &\leq |g|_{*,T \wedge \tau} \int_0^t e^{-m(t-s)}d\lambda(s) = |g|_{*,T \wedge \tau} \frac{1 - e^{-tm}}{m} \leq \frac{|g|_{*,T \wedge \tau}}{m}. \end{aligned} \quad (7.127)$$

Next we estimate the third term in the right hand side of (7.126). Since  $f$  is non-negative on  $[0, T \wedge \tau]$ ,  $de^F$  is a positive measure on  $(0, T \wedge \tau]$ . Hence for  $t \in [0, T \wedge \tau]$

$$\left| e^{-F(t)} \int_0^t M(s)de^F(s) \right| \leq |M|_{*,T \wedge \tau} e^{-F(t)} \int_0^t de^F(s) \leq |M|_{*,T \wedge \tau}. \quad (7.128)$$

Finally, the last term in the right hand side of (7.126) for any  $t \in [0, \tau \wedge T]$  can be bounded as

$$\left| e^{-F(t)}(z(0) - M(0)) \right| \leq (|z(0)| + |M(0)|)e^{-F(t)} \leq (|z(0)| + |M(0)|)e^{-mt}. \quad (7.129)$$

Using (7.127), (7.128) and (7.129) in (7.126) completes the proof of the lemma.  $\square$

## CHAPTER 8

### Conclusion

We have studied the topics of correlation mining and distributed load balancing in this dissertation. In conclusion, we will now summarize some of the main ideas that we have encountered and suggest future research directions.

#### 8.1 Correlation mining : summary and future directions

In the first part of this dissertation, we focused on the exploratory problem of detecting *bimodules*, which, recall, are feature set pairs from two data types that have significant aggregate cross-correlations. Since bimodules are based on groupwise association, they may provide more insights than standard pairwise analysis in correlation mining.

We introduced the Bimodule Search Procedure (BSP) in Chapter 3. Rather than relying on an underlying generative model, BSP makes use of iterative hypothesis-testing to identify *stable* bimodules, which satisfy a natural stability condition. The false discovery threshold  $\alpha \in (0, 1)$  is the only free parameter of BSP. Efficient approximation of the p-values used for iterative testing allowed BSP to run on large datasets.

At the population level, stable bimodules can be characterized in terms of the connected components of the population cross-correlation network. At the sample level, stable bimodules depend on both cross-correlations and intra-correlations, which are not part of the cross-correlation network. Nevertheless, the network perspective provides insights in both the simulation study and the real data analysis.

In Section 4.1, using a complex, network-based simulation study, we found that BSP was able to recover most true bimodules with significant cross-correlation strength, while simultaneously controlling the false discovery of edges having network-level importance. Among true bimodules with similar cross-correlation strengths, those with lower intra-correlations were more likely to be

recovered than those with higher intra-correlations, reflecting the incorporation of intra-correlations in the calculation of p-values; the effects of intra-correlations were most pronounced when the cross-correlation strength was moderate.

When applied to eQTL data in Section 4.2, BSP bimodules identified both local and distal effects, capturing half of the eQTLs found by standard cis-analysis and most of the eQTLs found by standard trans-analysis. Further, a substantial proportion of bimodules contained SNP-gene pairs that were important at the network level but not deemed significant under pairwise trans-analysis.

Finally, here are some future research directions that may address limitations of our work.

### 8.1.1 Moving beyond connected components in bimodules

At root, the discovery of bimodules by BSP and CONDOR (see Section 2.2.2) is driven by the presence or absence of correlations between features of different types. A key issue for these, and related, methods is how they behave with increasing sample size. In general, increasing sample size will yield greater power to detect cross-correlations, and therefore one expects the sizes of bimodule to increase. While this is often a desirable outcome, in applications where non-zero cross-correlations (possibly of small size) are the norm, this increased power may yield very large bimodules with little interpretive value. Evidence of this phenomena is found in our simulation study in Section 4.1 where, due to the presence of bridge-edges that connect different true bimodules, increasing the sample size from  $n = 200$  to  $n = 600$  yields larger BSP bimodules, which often contain multiple true bimodules (4.1.4). This may well reflect the underlying biology of genetic regulation: the omni-genic hypothesis of Boyle et al. (21) suggests that a substantial portion of the gene-SNP cross-correlation network might be connected at the population level.

An obvious way to address “super connectivity” of the cross-correlation network is to change the definition of bimodule to account for the magnitude of cross-correlations, rather than their mere presence or absence. Incorporating a more stringent definition of connectivity in BSP would require modifying the permutation null distribution and addressing the theory and computation behind such a modification, both of which may be future research directions.

### 8.1.2 False discovery guarantees in iterative testing

BSP is based on the framework of iterative-testing (IT) that has also been employed in other applications (98, 17, 18). The IT framework is a flexible tool that can translate a search procedure at the population level into one that works at the sample level. Indeed, as we saw in Chapter 3, at an high level, BSP is simply the population search procedure described at the end of Section 3.2 augmented by performing multiple hypothesis tests at each step. Although multiple testing correction (12) is used to control the false discoveries at every step, it is not clear how this ensures control over the false discoveries (e.g. features  $(s, t) \in (A, B)$  such that  $\rho(s, t) = 0$ ) among the detected bimodules. Presently, this lack of theoretical guarantee is a general limitation of the IT framework. From both a theoretical and practical stand point, it would be interesting to control the number of false discoveries that can emerge after several iterations of the multiple-testing step in such procedures. Such a result may be important for the challenge of establishing procedures that provide accurate false discovery control while making simultaneous inference over all sub-matrices of the covariance or precision matrix of interest (27).

## 8.2 Distributed load balancing : summary and future directions

Motivated by distributed load balancing, in the second part of this dissertation, we studied the balls and bins problem in discrete time and the Supermarket model in continuous time. Routing in these models is done using the power-of- $d$  scheme, where each incoming ball (or job) selects a random sample of  $d$  out of  $n$  bins (servers) and enters into the least occupied of the sampled bins (queues). The expression ‘power of choice’ is then used to denote the phenomenon that taking  $d = 2$  instead of  $d = 1$  causes a substantial improvement in the load balancing performance. The asymptotic behavior of these models for a fixed  $d \in \mathbb{N}$  as  $n \rightarrow \infty$  is well studied in the literature.

Motivated by recent developments in the queuing theory, our aim in this dissertation was to study the asymptotic properties of the power-of- $d$  models as  $d = d_n$  is allowed to increase with  $n$ . In particular, in Chapter 6, we show that the maximum bin size for the balls and bins problem when all the  $n$  balls are inserted is still  $\frac{\log \log n}{\log d_n} + O(1)$  with high probability if  $d_n = O(\log n)$ . In particular, this says that if  $d_n = (\log n)^{1/k}$  for some  $k \geq 1$ , then the maximum bin size will remain bounded as  $n \rightarrow \infty$ .

Our main results are in Chapter 7, where we address the limiting behavior of the Supermarket model with arrival rate  $n\lambda_n$  and  $n$  unit rate servers, as the number of choices  $d = d_n$  increases with  $n$ . We first establish a functional law of large numbers for the system as  $\lambda_n \rightarrow \lambda$  and  $(n, d_n) \rightarrow \infty$ , which shows that the scaled process converges in probability to a deterministic trajectory in time, called the fluid limit. The fluid limit is universal, in that it only depends on  $\lambda$  and is independent of the rate at which  $d_n \rightarrow \infty$ . Next, we establish functional central limit theorems for  $O(1/\sqrt{n})$  deviations of the scaled system around its ‘near-equilibrium’ point, which is the configuration at which the drift term in the semi-martingale representation vanishes. This near-equilibrium point which depends on  $(\lambda_n, d_n)$  acts almost like a fixed point of the system. For suitable choices of  $\lambda_n \rightarrow 1$ , under the three regimes (a)  $d_n/\sqrt{n} \rightarrow 0$ , (b)  $d_n/\sqrt{n} \rightarrow (0, \infty)$ , and (c)  $d_n/\sqrt{n} \rightarrow \infty$  we show functional central limit theorems for  $O(1/\sqrt{n})$  deviations of the system around the near equilibrium point.

Now we describe some research direction to pursue from here. First the central limit theorems described above for the Supermarket model only address convergence for bounded time intervals. The interchange of time and space limits and the problem of proving convergence for the corresponding equilibrium distributions remains open. Perhaps techniques used in Braverman (23) for the case  $d_n = n$  and  $\sqrt{n}(1 - \lambda_n) \rightarrow \beta$  could be extended to our setup. Additionally, the concentration analysis of (24) and attractive properties of the near-equilibrium point might also be helpful.

From a practical perspective, the limit theorems that we have proved for the Supermarket model may elucidate the relationship between  $(d_n, \lambda_n, n)$  required for good system performance, and may become a part of the arsenal of tools that can be used to design large scale data centers and other distributed systems. On the other hand, from a theoretical perspective, these results have been an exercise in using tools from stochastic calculus to study the behavior of complicated discrete-event systems. One may pursue this direction further and derive limit theorems for related models like the Supermarket model with memory (116, 76) or the power-of-choice on growing trees (43, 81).

## APPENDIX A: GENE ONTOLOGY ENRICHMENT OF BIMODULES

As explained in Section 4.2.4.3, the Gene Ontology (GO) database contains a curated collection of gene sets that are known to be associated with different biological functions. For each of the 145 BSP bimodules having a gene set  $B$  with 8 or more elements, we assessed the enrichment of  $B$  in 6463 GO gene sets of size more than 10, representing biological processes; 18 out of the 145 BSP gene sets, and 1 out of the 5 CONDOR gene sets that we considered had significant overlap with Gene Ontology (GO) categories. Among the 40 GO terms detected by the CONDOR, 27 terms were also found among the 135 terms detected by BSP. The GO terms that were discovered for the two methods did not seem specific to thyroid, the tissue under investigation.

Complete list of the GO terms for the two methods is as follows. Significant GO terms for BSP:

	GO.ID	Term	bimod.
1	GO:0060333	interferon-gamma-mediated signaling path...	1
2	GO:0002478	antigen processing and presentation of e...	1
3	GO:0019884	antigen processing and presentation of e...	1
4	GO:0048002	antigen processing and presentation of p...	1
5	GO:0019882	antigen processing and presentation	1
6	GO:0071346	cellular response to interferon-gamma	1
7	GO:0034341	response to interferon-gamma	1
8	GO:0019886	antigen processing and presentation of e...	1
9	GO:0002495	antigen processing and presentation of p...	1
10	GO:0002504	antigen processing and presentation of p...	1
11	GO:0045087	innate immune response	1
12	GO:0050776	regulation of immune response	1
13	GO:0006952	defense response	1
14	GO:0031295	T cell costimulation	1
15	GO:0031294	lymphocyte costimulation	1
16	GO:0050852	T cell receptor signaling pathway	1
17	GO:0002768	immune response-regulating cell surface ...	1
18	GO:0002764	immune response-regulating signaling pat...	1



19	GO:0050851	antigen receptor-mediated signaling path...	1
20	GO:0002682	regulation of immune system process	1
21	GO:0022409	positive regulation of cell-cell adhesio...	1
22	GO:0002253	activation of immune response	1
23	GO:0002429	immune response-activating cell surface ...	1
24	GO:0006950	response to stress	1
25	GO:0006955	immune response	1
26	GO:0019221	cytokine-mediated signaling pathway	1
27	GO:0002757	immune response-activating signal transd...	1
28	GO:0050870	positive regulation of T cell activation	1
29	GO:0002479	antigen processing and presentation of e...	1
30	GO:1903039	positive regulation of leukocyte cell-ce...	1
31	GO:0042590	antigen processing and presentation of e...	1
32	GO:0045806	negative regulation of endocytosis	2
33	GO:0050911	detection of chemical stimulus involved ...	3
34	GO:0007608	sensory perception of smell	3
35	GO:0050907	detection of chemical stimulus involved ...	3
36	GO:0009593	detection of chemical stimulus	3
37	GO:0007606	sensory perception of chemical stimulus	3
38	GO:0035459	cargo loading into vesicle	3
39	GO:0050906	detection of stimulus involved in sensor...	3
40	GO:0000038	very long-chain fatty acid metabolic pro...	4
41	GO:0006732	coenzyme metabolic process	4
42	GO:0006417	regulation of translation	5
43	GO:0034248	regulation of cellular amide metabolic p...	5
44	GO:0010608	posttranscriptional regulation of gene e...	5
45	GO:0046597	negative regulation of viral entry into ...	6
46	GO:0035455	response to interferon-alpha	6
47	GO:0035456	response to interferon-beta	6

48	GO:0046596	regulation of viral entry into host cell	6
49	GO:0045071	negative regulation of viral genome repl...	6
50	GO:1903901	negative regulation of viral life cycle	6
51	GO:0060337	type I interferon signaling pathway	6
52	GO:0071357	cellular response to type I interferon	6
53	GO:0034340	response to type I interferon	6
54	GO:0045069	regulation of viral genome replication	6
55	GO:0048525	negative regulation of viral process	6
56	GO:0019079	viral genome replication	6
57	GO:0046718	viral entry into host cell	6
58	GO:1903900	regulation of viral life cycle	6
59	GO:0030260	entry into host cell	6
60	GO:0044409	entry into host	6
61	GO:0051806	entry into cell of other organism involv...	6
62	GO:0051828	entry into other organism involved in sy...	6
63	GO:0043901	negative regulation of multi-organism pr...	6
64	GO:0034341	response to interferon-gamma	6
65	GO:0050792	regulation of viral process	6
66	GO:0051607	defense response to virus	6
67	GO:0051701	interaction with host	6
68	GO:0043903	regulation of symbiosis, encompassing mu...	6
69	GO:0009615	response to virus	6
70	GO:0051225	spindle assembly	7
71	GO:0007030	Golgi organization	7
72	GO:0007051	spindle organization	7
73	GO:0010256	endomembrane system organization	7
74	GO:0000226	microtubule cytoskeleton organization	7
75	GO:0007017	microtubule-based process	7
76	GO:0070925	organelle assembly	7

77	GO:0007010	cytoskeleton organization	7
78	GO:0007156	homophilic cell adhesion via plasma memb...	8
79	GO:0098742	cell-cell adhesion via plasma-membrane a...	8
80	GO:0098609	cell-cell adhesion	8
81	GO:0007155	cell adhesion	8
82	GO:0022610	biological adhesion	8
83	GO:0007416	synapse assembly	8
84	GO:0007267	cell-cell signaling	8
85	GO:0006355	regulation of transcription, DNA-templat...	9
86	GO:1903506	regulation of nucleic acid-templated tra...	9
87	GO:2001141	regulation of RNA biosynthetic process	9
88	GO:0006351	transcription, DNA-templated	9
89	GO:0097659	nucleic acid-templated transcription	9
90	GO:0032774	RNA biosynthetic process	9
91	GO:0051252	regulation of RNA metabolic process	9
92	GO:2000112	regulation of cellular macromolecule bio...	9
93	GO:0010556	regulation of macromolecule biosynthetic...	9
94	GO:0019219	regulation of nucleobase-containing comp...	9
95	GO:0031326	regulation of cellular biosynthetic proc...	9
96	GO:0034654	nucleobase-containing compound biosynthe...	9
97	GO:0009889	regulation of biosynthetic process	9
98	GO:0018130	heterocycle biosynthetic process	9
99	GO:0019438	aromatic compound biosynthetic process	9
100	GO:0010468	regulation of gene expression	9
101	GO:1901362	organic cyclic compound biosynthetic pro...	9
102	GO:0016070	RNA metabolic process	9
103	GO:0001580	detection of chemical stimulus involved ...	10
104	GO:0050912	detection of chemical stimulus involved ...	10
105	GO:0050913	sensory perception of bitter taste	10

106	GO:0050909	sensory perception of taste	10
107	GO:0050907	detection of chemical stimulus involved ...	10
108	GO:0009593	detection of chemical stimulus	10
109	GO:0007606	sensory perception of chemical stimulus	10
110	GO:0050906	detection of stimulus involved in sensor...	10
111	GO:0007600	sensory perception	10
112	GO:0051606	detection of stimulus	10
113	GO:0050877	nervous system process	10
114	GO:0003008	system process	10
115	GO:0007186	G-protein coupled receptor signaling pat...	10
116	GO:0006355	regulation of transcription, DNA-templat...	11
117	GO:1903506	regulation of nucleic acid-templated tra...	11
118	GO:2001141	regulation of RNA biosynthetic process	11
119	GO:0006351	transcription, DNA-templated	11
120	GO:0097659	nucleic acid-templated transcription	11
121	GO:0032774	RNA biosynthetic process	11
122	GO:0051252	regulation of RNA metabolic process	11
123	GO:2000112	regulation of cellular macromolecule bio...	11
124	GO:0010556	regulation of macromolecule biosynthetic...	11
125	GO:0019219	regulation of nucleobase-containing comp...	11
126	GO:0031326	regulation of cellular biosynthetic proc...	11
127	GO:0034654	nucleobase-containing compound biosynthe...	11
128	GO:0009889	regulation of biosynthetic process	11
129	GO:0018130	heterocycle biosynthetic process	11
130	GO:0019438	aromatic compound biosynthetic process	11
131	GO:0010468	regulation of gene expression	11
132	GO:1901362	organic cyclic compound biosynthetic pro...	11
133	GO:0016070	RNA metabolic process	11
134	GO:1901685	glutathione derivative metabolic process	12

135	GO:1901687	glutathione derivative biosynthetic proc...	12
136	GO:0006749	glutathione metabolic process	12
137	GO:0042178	xenobiotic catabolic process	12
138	GO:0042537	benzene-containing compound metabolic pr...	12
139	GO:0006575	cellular modified amino acid metabolic p...	12
140	GO:0044272	sulfur compound biosynthetic process	12
141	GO:0046854	phosphatidylinositol phosphorylation	13
142	GO:0046834	lipid phosphorylation	13
143	GO:0048015	phosphatidylinositol-mediated signaling	13
144	GO:0048017	inositol lipid-mediated signaling	13
145	GO:0006882	cellular zinc ion homeostasis	14
146	GO:0055069	zinc ion homeostasis	14
147	GO:0010273	detoxification of copper ion	14
148	GO:1990169	stress response to copper ion	14
149	GO:0061687	detoxification of inorganic compound	14
150	GO:0097501	stress response to metal ion	14
151	GO:0071294	cellular response to zinc ion	14
152	GO:0071280	cellular response to copper ion	14
153	GO:0046916	cellular transition metal ion homeostasi...	14
154	GO:0071276	cellular response to cadmium ion	14
155	GO:0046688	response to copper ion	14
156	GO:0055076	transition metal ion homeostasis	14
157	GO:0072488	ammonium transmembrane transport	15
158	GO:0006089	lactate metabolic process	16
159	GO:0006882	cellular zinc ion homeostasis	17
160	GO:0055069	zinc ion homeostasis	17
161	GO:0006882	cellular zinc ion homeostasis	18
162	GO:0055069	zinc ion homeostasis	18

---

Significant GO terms for CONDOR:

	GO.ID	Term	bimod
1	GO:0050852	T cell receptor signaling pathway	1
2	GO:0050851	antigen receptor-mediated signaling path...	1
3	GO:0006355	regulation of transcription, DNA-templat...	1
4	GO:1903506	regulation of nucleic acid-templated tra...	1
5	GO:2001141	regulation of RNA biosynthetic process	1
6	GO:0060333	interferon-gamma-mediated signaling path...	2
7	GO:0002478	antigen processing and presentation of e...	3
8	GO:0019884	antigen processing and presentation of e...	3
9	GO:0048002	antigen processing and presentation of p...	3
10	GO:0019886	antigen processing and presentation of e...	3
11	GO:0002495	antigen processing and presentation of p...	3
12	GO:0002504	antigen processing and presentation of p...	3
13	GO:0019882	antigen processing and presentation	3
14	GO:0031295	T cell costimulation	3
15	GO:0031294	lymphocyte costimulation	3
16	GO:0060333	interferon-gamma-mediated signaling path...	3
17	GO:0050852	T cell receptor signaling pathway	3
18	GO:0050870	positive regulation of T cell activation	3
19	GO:1903039	positive regulation of leukocyte cell-ce...	3
20	GO:0050778	positive regulation of immune response	3
21	GO:0002253	activation of immune response	3
22	GO:0050851	antigen receptor-mediated signaling path...	3
23	GO:0022409	positive regulation of cell-cell adhesio...	3
24	GO:0071346	cellular response to interferon-gamma	3
25	GO:0051251	positive regulation of lymphocyte activa...	3
26	GO:1903037	regulation of leukocyte cell-cell adhesi...	3
27	GO:0034341	response to interferon-gamma	3
28	GO:0002696	positive regulation of leukocyte activat...	3

29	GO:0050863	regulation of T cell activation	3
30	GO:0050867	positive regulation of cell activation	3
31	GO:0007159	leukocyte cell-cell adhesion	3
32	GO:0050776	regulation of immune response	3
33	GO:0002429	immune response-activating cell surface ...	3
34	GO:0002684	positive regulation of immune system pro...	3
35	GO:0006955	immune response	3
36	GO:0022407	regulation of cell-cell adhesion	3
37	GO:0045087	innate immune response	3
38	GO:0002768	immune response-regulating cell surface ...	3
39	GO:0045785	positive regulation of cell adhesion	3
40	GO:0002455	humoral immune response mediated by circ...	3
41	GO:0051249	regulation of lymphocyte activation	3
42	GO:0019221	cytokine-mediated signaling pathway	3
43	GO:0042110	T cell activation	3

---

## BIBLIOGRAPHY

- [1] Albert, F. W. and Kruglyak, L. (2015). The role of regulatory variation in complex traits and disease. *Nature Reviews Genetics*, 16(4):197.
- [2] Alexa, A. and Rahnenführer, J. (2009). Gene set enrichment analysis with topgo. *Bioconductor Improv*, 27.
- [3] Alexa, A. and Rahnenfuhrer, J. (2018). *topGO: Enrichment Analysis for Gene Ontology*. R package version 2.34.0.
- [4] Altman, E., Ayesta, U., and Prabhu, B. J. (2011). Load balancing in processor sharing systems. *Telecommunication Systems*, 47(1-2):35–48.
- [5] Avin, C. and Krishnamachari, B. (2006). The power of choice in random walks: An empirical study. In *Proceedings of the 9th ACM international symposium on Modeling analysis and simulation of wireless and mobile systems*, pages 219–228.
- [6] Azar, Y., Broder, A. Z., Karlin, A. R., and Upfal, E. (1994). Balanced allocations. In *Proceedings of the twenty-sixth annual ACM symposium on theory of computing*, pages 593–602.
- [7] Banerjee, S. and Mukherjee, D. (2018). Join-the-shortest queue diffusion limit in Halfin-Whitt regime: Sensitivity on the heavy-traffic parameter. *arXiv preprint arXiv:1809.01739*.
- [8] Banerjee, S. and Mukherjee, D. (2019). Join-the-shortest queue diffusion limit in Halfin-Whitt regime: Tail asymptotics and scaling of extrema. *The Annals of Applied Probability*, 29(2):1262–1309.
- [9] Barber, M. J. (2007). Modularity and community detection in bipartite networks. *Physical Review E*, 76(6):066102.
- [10] Beckett, S. J. (2016). Improved community detection in weighted bipartite networks. *Royal Society open science*, 3(1):140536.
- [11] Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300.
- [12] Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, pages 1165–1188.
- [13] Berenbrink, P., Czumaj, A., Steger, A., and Vöcking, B. (2006). Balanced allocations: The heavily loaded case. *SIAM Journal on Computing*, 35(6):1350–1385.
- [14] Berenbrink, P., Friedetzky, T., Kling, P., Mallmann-Trenn, F., Nagel, L., and Wastell, C. (2016). Self-stabilizing balls & bins in batches: The power of leaky bins. In *Proceedings of the 2016 ACM Symposium on Principles of Distributed Computing*, pages 83–92.
- [15] Berenbrink, P., Khodamoradi, K., Sauerwald, T., and Stauffer, A. (2013). Balls-into-bins with nearly optimal load distribution. In *Proceedings of the twenty-fifth annual ACM symposium on Parallelism in algorithms and architectures*, pages 326–335.



- [16] Billingsley, P. (1995). Probability and measure. 1995. *John Wiley&Sons, New York*.
- [17] Bodwin, K., Chakraborty, S., Zhang, K., and Nobel, A. B. (2017). Latent association mining in binary data. *arXiv preprint arXiv:1711.10427*.
- [18] Bodwin, K., Zhang, K., Nobel, A., et al. (2018). A testing based approach to the discovery of differentially correlated variable sets. *The Annals of Applied Statistics*, 12(2):1180–1203.
- [19] Bollobás, B. (2001). The evolution of random graphs—the giant component. In *Random graphs*, volume 184, pages 130–59. Cambridge university press Cambridge.
- [20] Botstein, D., Cherry, J. M., Ashburner, M., Ball, C., Blake, J., Butler, H., Davis, A., Dolinski, K., Dwight, S., Eppig, J., et al. (2000). Gene ontology: tool for the unification of biology. *Nat genet*, 25(1):25–9.
- [21] Boyle, E. A., Li, Y. I., and Pritchard, J. K. (2017). An expanded view of complex traits: from polygenic to omnigenic. *Cell*, 169(7):1177–1186.
- [22] Bramson, M., Lu, Y., and Prabhakar, B. (2012). Asymptotic independence of queues under randomized load balancing. *Queueing Systems*, 71(3):247–292.
- [23] Braverman, A. (2020). Steady-state analysis of the join-the-shortest-queue model in the halfin–whitt regime. *Mathematics of Operations Research*, 45(3):1069–1103.
- [24] Brightwell, G., Fairthorne, M., and Luczak, M. J. (2018). The supermarket model with bounded queue lengths in equilibrium. *Journal of Statistical Physics*, 173(3-4):1149–1194.
- [25] Brightwell, G. and Luczak, M. (2012). The supermarket model with arrival rate tending to one. *arXiv preprint arXiv:1201.5523*.
- [26] Budhiraja, A. and Friedlander, E. (2019). Diffusion approximations for load balancing mechanisms in cloud storage systems. *Advances in Applied Probability*, 51(1):41–86.
- [27] Cai, T. T. (2017). Global testing and large-scale multiple testing for high-dimensional covariance structures. *Annual Review of Statistics and Its Application*, 4:423–446.
- [28] Cai, T. T. and Liu, W. (2016). Large-scale multiple testing of correlations. *Journal of the American Statistical Association*, 111(513):229–240.
- [29] Cardellini, V., Casalicchio, E., Colajanni, M., and Yu, P. S. (2002). The state of the art in locally distributed web-server systems. *ACM Computing Surveys (CSUR)*, 34(2):263–311.
- [30] Chen, X., Han, L., and Carbonell, J. (2012a). Structured sparse canonical correlation analysis. In *Artificial intelligence and statistics*, pages 199–207. PMLR.
- [31] Chen, X., Shi, X., Xu, X., Wang, Z., Mills, R., Lee, C., and Xu, J. (2012b). A two-graph guided multi-task lasso approach for eqtl mapping. *Journal of Machine Learning Research*, 22:208–217.
- [32] Cheng, W., Shi, Y., Zhang, X., and Wang, W. (2015). Fast and robust group-wise eqtl mapping using sparse graphical models. *BMC bioinformatics*, 16(1):2.
- [33] Cheng, W., Shi, Y., Zhang, X., and Wang, W. (2016). Sparse regression models for unraveling group and individual associations in eqtl mapping. *BMC bioinformatics*, 17(1):136.

- [34] Cheng, W., Zhang, X., Wu, Y., Yin, X., Li, J., Heckerman, D., and Wang, W. (2012). Inferring novel associations between SNP sets and gene sets in eqtl study using sparse graphical model. In *Proceedings of the ACM Conference on Bioinformatics, Computational Biology and Biomedicine*, pages 466–473. ACM.
- [35] Cho, Y. J., Wang, J., and Joshi, G. (2020). Client selection in federated learning: Convergence analysis and power-of-choice selection strategies. *arXiv preprint arXiv:2010.01243*.
- [36] Chu, J.-h., Weiss, S. T., Carey, V. J., and Raby, B. A. (2009). A graphical model approach for inferring large-scale networks integrating gene expression and genetic polymorphism. *BMC systems biology*, 3(1):55.
- [37] Consortium, G. et al. (2017). Genetic effects on gene expression across human tissues. *Nature*, 550(7675):204.
- [38] Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). *Introduction to algorithms*. MIT press.
- [39] Costa, A. and Hansen, P. (2014). A locally optimal hierarchical divisive heuristic for bipartite modularity maximization. *Optimization letters*, 8(3):903–917.
- [40] Czumaj, A. and Stemmann, V. (1997). Randomized allocation processes. In *Proceedings 38th Annual Symposium on Foundations of Computer Science*, pages 194–203. IEEE.
- [41] Dewaskar, M., Palowitch, J., He, M., Love, M. I., and Nobel, A. (2020). Finding stable groups of cross-correlated features in multi-view data. *arXiv preprint arXiv:2009.05079*.
- [42] Dolédec, S. and Chessel, D. (1994). Co-inertia analysis: an alternative method for studying species–environment relationships. *Freshwater biology*, 31(3):277–294.
- [43] D’Souza, R. M., Krapivsky, P. L., and Moore, C. (2007). The power of choice in growing trees. *The European Physical Journal B*, 59(4):535–543.
- [44] Eschenfeldt, P. and Gamarnik, D. (2016). Supermarket queueing system in the heavy traffic regime. Short queue dynamics. *arXiv preprint arXiv:1610.03522*.
- [45] Eschenfeldt, P. and Gamarnik, D. (2018). Join the shortest queue with many servers. The heavy-traffic asymptotics. *Mathematics of Operations Research*, 43(3):867–886.
- [46] Ethier, S. N. and Kurtz, T. G. (2009). *Markov Processes: Characterization and Convergence*, volume 282. John Wiley & Sons.
- [47] Fagny, M., Paulson, J. N., Kuijjer, M. L., Sonawane, A. R., Chen, C.-Y., Lopes-Ramos, C. M., Glass, K., Quackenbush, J., and Platig, J. (2017). Exploring regulation in tissues with eqtl networks. *Proceedings of the National Academy of Sciences*, 114(37):E7841–E7850.
- [48] Gamarnik, D., Tsitsiklis, J. N., and Zubeldia, M. (2016). Delay, memory, and messaging tradeoffs in distributed service systems. *ACM SIGMETRICS Performance Evaluation Review*, 44(1):1–12.
- [49] Gamarnik, D., Tsitsiklis, J. N., and Zubeldia, M. (2018). Delay, memory, and messaging tradeoffs in distributed service systems. *Stochastic Systems*, 8(1):45–74.

- [50] Gene Ontology Consortium (2014). Gene ontology consortium: going forward. *Nucleic acids research*, 43(D1):D1049–D1056.
- [51] Georgakopoulos, A., Haslegrave, J., Sauerwald, T., and Sylvester, J. (2019). The power of two choices for random walks. *arXiv preprint arXiv:1911.05170*.
- [52] Graham, C. (2000). Chaoticity on path space for a queueing network with selection of the shortest queue among several. *Journal of Applied Probability*, 37(1):198–211.
- [53] Graham, C. (2005). Functional central limit theorems for a large network in which customers join the shortest of several queues. *Probability theory and related fields*, 131(1):97–120.
- [54] GTEx Consortium (2017). Genetic effects on gene expression across human tissues. *Nature*, 550(7675):204.
- [55] Gupta, V., Balter, M. H., Sigman, K., and Whitt, W. (2007). Analysis of join-the-shortest-queue routing for web server farms. *Performance Evaluation*, 64(9-12):1062–1081.
- [56] Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media.
- [57] Hero, A. and Rajaratnam, B. (2011a). Hub discovery in partial correlation graphical models. *arXiv preprint arXiv:1109.6846*.
- [58] Hero, A. and Rajaratnam, B. (2011b). Large-scale correlation screening. *Journal of the American Statistical Association*, 106(496):1540–1552.
- [59] Hero, A. O. and Rajaratnam, B. (2015). Foundational principles for large-scale inference: Illustrations through correlation mining. *Proceedings of the IEEE*, 104(1):93–110.
- [60] Hotelling, H. (1936). Relations between two sets of variates. *Biometrika*, 28(3/4):321–377.
- [61] (<https://math.stackexchange.com/users/354840/kajelad>), K. A continous extension of the group of permutation matrices. Mathematics Stack Exchange. URL:<https://math.stackexchange.com/q/3778335> (version: 2020-08-03).
- [62] Hunt, P. and Kurtz, T. (1994). Large loss networks. *Stochastic Processes and their Applications*, 53(2):363–378.
- [63] Karatzas, I. and Shreve, S. E. (1998). *Brownian Motion and Stochastic Calculus*. Springer-Verlag, New York.
- [64] Kenthapadi, K. and Panigrahy, R. (2006). Balanced allocation on graphs. In *SODA*, volume 6, pages 434–443.
- [65] Kolberg, L., Kerimov, N., Peterson, H., and Alasoo, K. (2020). Co-expression analysis reveals interpretable gene modules controlled by *trans*-acting genetic variants. *eLife*, 9:e58705.
- [66] Kurtz, T. G. (1978). Strong approximation theorems for density dependent Markov chains. *Stochastic Processes and their Applications*, 6(3):223–240.
- [67] Kurtz, T. G. (1981). *Approximation of Population Processes*, volume 36. SIAM.
- [68] Lahat, D., Adali, T., and Jutten, C. (2015). Multimodal data fusion: An overview of methods, challenges, and prospects. *Proceedings of the IEEE*, 103(9):1449–1477.

- [69] Lenzen, C. and Wattenhofer, R. (2016). Tight bounds for parallel randomized load balancing. *Distributed Computing*, 29(2):127–142.
- [70] Li, B., Ramamoorthy, A., and Srikant, R. (2016). Mean-field-analysis of coding versus replication in cloud storage systems. In *IEEE INFOCOM 2016 - The 35th Annual IEEE International Conference on Computer Communications*, pages 1–9.
- [71] Lin, D., Zhang, J., Li, J., Calhoun, V. D., Deng, H.-W., and Wang, Y.-P. (2013). Group sparse canonical correlation analysis for genomic data integration. *BMC bioinformatics*, 14(1):245.
- [72] Liu, X. and Murata, T. (2010). An efficient algorithm for optimizing bipartite modularity in bipartite networks. *Journal of Advanced Computational Intelligence and Intelligent Informatics*, 14(4):408–415.
- [73] Liu, X. and Ying, L. (2019). A simple steady-state analysis of load balancing algorithms in the sub-halfin-whitt regime. *ACM SIGMETRICS Performance Evaluation Review*, 46(2):15–17.
- [74] Liu, X. and Ying, L. (2020). Steady-state analysis of load-balancing algorithms in the sub-halfin-whitt regime. *Journal of Applied Probability*, 57(2):578–596.
- [75] Luczak, M. and McDiarmid, C. (2007). Asymptotic distributions and chaos for the supermarket model. *Electronic Journal of Probability*, 12:75–99.
- [76] Luczak, M., Norris, J., et al. (2013). Averaging over fast variables in the fluid limit for markov chains: application to the supermarket model with memory. *The Annals of Applied Probability*, 23(3):957–986.
- [77] Luczak, M. J. and McDiarmid, C. (2005). On the power of two choices: Balls and bins in continuous time. *Ann. Appl. Probab.*, 15(3):1733–1764.
- [78] Luczak, M. J. and McDiarmid, C. (2006). On the maximum queue length in the supermarket model. *The Annals of Probability*, 34(2):493–527.
- [79] Luczak, M. J. and Norris, J. (2005). Strong approximation for the supermarket model. *The Annals of Applied Probability*, 15(3):2038–2061.
- [80] Maguluri, S. T., Srikant, R., and Ying, L. (2012). Stochastic models of load balancing and scheduling in cloud computing clusters. In *2012 Proceedings IEEE Infocom*, pages 702–710. IEEE.
- [81] Malyshkin, Y., Paquette, E., et al. (2014). The power of choice combined with preferential attachment. *Electronic Communications in Probability*, 19.
- [82] Martin, J. and Suhov, Y. M. (1999). Fast Jackson networks. *The Annals of Applied Probability*, 9(3):854–870.
- [83] McCabe, S. D., Lin, D.-Y., and Love, M. I. (2019). Consistency and overfitting of multi-omics methods on experimental data. *Brief Bioinform.*
- [84] McDiarmid, C. (1998). Concentration. In *Probabilistic methods for algorithmic discrete mathematics*, pages 195–248. Springer.
- [85] McIntosh, A., Bookstein, F., Haxby, J. V., and Grady, C. (1996). Spatial pattern analysis of functional brain images using partial least squares. *Neuroimage*, 3(3):143–157.

- [86] McShane, E. J. and Botts, T. A. (2013). *Real Analysis*. Courier Corporation.
- [87] Meckes, E. (2014). Concentration of measure and the compact classical matrix groups.
- [88] Meng, C., Zeleznik, O. A., Thallinger, G. G., Kuster, B., Gholami, A. M., and Culhane, A. C. (2016). Dimension reduction techniques for the integrative analysis of multi-omics data. *Briefings in bioinformatics*, 17(4):628–641.
- [89] Mitzenmacher, M. (2001). The power of two choices in randomized load balancing. *IEEE Transactions on Parallel and Distributed Systems*, 12(10):1094–1104.
- [90] Muirhead, R. J. (2009). *Aspects of multivariate statistical theory*, volume 197. John Wiley & Sons.
- [91] Mukherjee, D., Borst, S. C., Van Leeuwen, J. S., and Whiting, P. A. (2018). Universality of power-of-d load balancing in many-server systems. *Stochastic Systems*, 8(4):265–292.
- [92] Mukhopadhyay, A. and Mazumdar, R. R. (2015). Analysis of randomized join-the-shortest-queue (JSQ) schemes in large heterogeneous processor-sharing systems. *IEEE Transactions on Control of Network Systems*, 3(2):116–126.
- [93] Nash, J. F. (1950). Equilibrium points in n-person games. *Proceedings of the National Academy of Sciences*, 36(1):48–49.
- [94] Nica, A. C. and Dermitzakis, E. T. (2013). Expression quantitative trait loci: present and future. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 368(1620):20120362.
- [95] O’neil, E. J., O’neil, P. E., and Weikum, G. (1993). The lru-k page replacement algorithm for database disk buffering. *Acm Sigmod Record*, 22(2):297–306.
- [96] Ongaro, D., Rumble, S. M., Stutsman, R., Ousterhout, J., and Rosenblum, M. (2011). Fast crash recovery in RAMCloud. In *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*, pages 29–41.
- [97] Ousterhout, K., Wendell, P., Zaharia, M., and Stoica, I. (2013). Sparrow: distributed, low latency scheduling. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 69–84.
- [98] Palowitch, J., Bhamidi, S., and Nobel, A. B. (2016). The continuous configuration model: A null for community detection on weighted networks. *arXiv preprint arXiv:1601.05630*.
- [99] Pan, C., Luo, J., Zhang, J., and Li, X. (2019). BiModule: biclique modularity strategy for identifying transcription factor and microRNA co-regulatory modules. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*.
- [100] Park, G. (2017). A generalization of multiple choice balls-into-bins: Tight bounds. *Algorithmica*, 77(4):1159–1193.
- [101] Parkhomenko, E., Tritchler, D., and Beyene, J. (2007). Genome-wide sparse canonical correlation of gene expression with genotypes. In *BMC proceedings*, volume 1, page S119. BioMed Central.

- [102] Parkhomenko, E., Tritchler, D., and Beyene, J. (2009). Sparse canonical correlation analysis with application to genomic data integration. *Statistical applications in genetics and molecular biology*, 8(1).
- [103] Patel, P. V., Gianoulis, T. A., Bjornson, R. D., Yip, K. Y., Engelman, D. M., and Gerstein, M. B. (2010). Analysis of membrane proteins in metagenomics: Networks of correlated environmental features and protein families. *Genome Research*, 20(7):960–971.
- [104] Peres, Y., Talwar, K., and Wieder, U. (2015). Graphical balanced allocations and the  $(1+\beta)$ -choice process. *Random Structures & Algorithms*, 47(4):760–775.
- [105] Pesantez-Cabrera, P. and Kalyanaraman, A. (2016). Detecting communities in biological bipartite networks. In *Proceedings of the 7th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, pages 98–107.
- [106] Platig, J. (2016). *condor: CComplex Network Description Of Regulators*. R package version 1.1.1.
- [107] Platig, J., Castaldi, P. J., DeMeo, D., and Quackenbush, J. (2016). Bipartite community structure of eqtls. *PLoS computational biology*, 12(9):e1005033.
- [108] Pourmiri, A. (2019). Balanced allocation on graphs: A random walk approach. *Random Structures & Algorithms*, 55(4):980–1009.
- [109] Pucher, B. M., Zeleznik, O. A., and Thallinger, G. G. (2019). Comparison and evaluation of integrative methods for the analysis of multilevel omics data: a study based on simulated and experimental cancer data. *Briefings in bioinformatics*, 20(2):671–681.
- [110] Raab, M. and Steger, A. (1998). “balls into bins”—a simple and tight analysis. In *International Workshop on Randomization and Approximation Techniques in Computer Science*, pages 159–170. Springer.
- [111] Rhee, S. Y., Wood, V., Dolinski, K., and Draghici, S. (2008). Use and misuse of the gene ontology annotations. *Nature Reviews Genetics*, 9(7):509–515.
- [112] Richa, A. W., Mitzenmacher, M., and Sitaraman, R. (2001). The power of two random choices: A survey of techniques and results. *Combinatorial Optimization*, 9:255–304.
- [113] Sankaran, K. and Holmes, S. P. (2019). Multitable methods for microbiome data integration. *Frontiers in genetics*, 10.
- [114] Shabalin, A. A. (2012). Matrix eqtl: ultra fast eqtl analysis via large matrix operations. *Bioinformatics*, 28(10):1353–1358.
- [115] Shabalin, A. A., Weigman, V. J., Perou, C. M., Nobel, A. B., et al. (2009). Finding large average submatrices in high dimensional data. *The Annals of Applied Statistics*, 3(3):985–1012.
- [116] Shah, D. and Prabhakar, B. (2002). The use of memory in randomized load balancing. In *Proceedings IEEE International Symposium on Information Theory*,, page 125. IEEE.
- [117] Stegle, O., Parts, L., Piipari, M., Winn, J., and Durbin, R. (2012). Using probabilistic estimation of expression residuals (peer) to obtain increased power and interpretability of gene expression analyses. *Nature protocols*, 7(3):500.

- [118] Talwar, K. and Wieder, U. (2014). Balanced allocations: A simple proof for the heavily loaded case. In *International Colloquium on Automata, Languages, and Programming*, pages 979–990. Springer.
- [119] Tian, L., Quitadamo, A., Lin, F., and Shi, X. (2014). Methods for population-based eqtl analysis in human genetics. *Tsinghua Science and Technology*, 19(6):624–634.
- [120] Tini, G., Marchetti, L., Priami, C., and Scott-Boyer, M.-P. (2019). Multi-omics integration—a comparison of unsupervised clustering methodologies. *Briefings in bioinformatics*, 20(4):1269–1279.
- [121] van der Boor, M., Borst, S. C., van Leeuwen, J. S., and Mukherjee, D. (2018). Scalable load balancing in networked systems: A survey of recent advances. *arXiv preprint arXiv:1806.05444*.
- [122] Van der Boor, M., Borst, S. C., Van Leeuwen, J. S., and Mukherjee, D. (2018). Scalable load balancing in networked systems: Universality properties and stochastic coupling methods. In *Proc. ICM*, volume 18. World Scientific.
- [123] van der Boor, M., Borst, S. C., van Leeuwen, J. S. H., and Mukherjee, D. (2018). Scalable load balancing in networked systems: A survey of recent advances.
- [124] Vvedenskaya, N. D., Dobrushin, R. L., and Karpelevich, F. I. (1996). Queueing system with selection of the shortest of two queues: An asymptotic approach. *Problemy Peredachi Informatsii*, 32(1):20–34.
- [125] Waaijenborg, S., de Witt Hamer, P. C. V., and Zwinderman, A. H. (2008). Quantifying the association between gene expressions and dna-markers by penalized canonical correlation analysis. *Statistical applications in genetics and molecular biology*, 7(1).
- [126] Westra, H.-J. and Franke, L. (2014). From genome to function by studying eqtls. *Biochimica et Biophysica Acta (BBA)-Molecular Basis of Disease*, 1842(10):1896–1902.
- [127] Wieder, U. et al. (2017). Hashing, load balancing and multiple choice. *Foundations and Trends® in Theoretical Computer Science*, 12(3–4):275–379.
- [128] Wilms, I. and Croux, C. (2015). Sparse canonical correlation analysis from a predictive point of view. *Biometrical Journal*, 57(5):834–851.
- [129] Witten, D. and Tibshirani, R. (2020). *PMA: Penalized Multivariate Analysis*. R package version 1.2.1.
- [130] Witten, D. M., Tibshirani, R., and Hastie, T. (2009). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534.
- [131] Wu, X., Liu, Q., and Jiang, R. (2009). Align human interactome with phenome to identify causative genes and networks underlying disease families. *Bioinformatics*, 25(1):98–104. Publisher: Oxford Academic.
- [132] Yin, J. and Li, H. (2011). A sparse conditional gaussian graphical model for analysis of genetical genomics data. *The annals of applied statistics*, 5(4):2630.

- [133] Yosida, K. (2012). *Functional Analysis*. Springer Science & Business Media.
- [134] Zhang, L. and Kim, S. (2014). Learning gene networks under snp perturbations using eqtl datasets. *PLoS computational biology*, 10(2):e1003420.
- [135] Zheng, X. (2015). A tutorial for the r/bioconductor package snprelate.
- [136] Zhou, Y.-H., Barry, W. T., and Wright, F. A. (2013). Empirical pathway analysis, without permutation. *Biostatistics*, 14(3):573–585.
- [137] Zhou, Y.-H., Gallins, P., and Wright, F. (2019). Marker-trait complete analysis. *bioRxiv*, page 836494.